# TARTU RIIKLIKU ÜLIKOOLI
# TOIMETISED

# 798

# STATISTICAL AND PROBABILISTIC MODELS

Matemaatika- ja mehaanikaalased tööd

Труды по математике и механике

TARTU 1988

# STATISTICAL AND PROBABILISTIC MODELS

Matemaatika- ja mehaanikaalased tööd

Труды по математике и механике

2 - 2

# UNBIASED AND $n^{-k}$-BIASED ESTIMATIONS OF ENTIRE RATIONAL FUNCTIONS OF MOMENTS.

## E.Tiit

### 1. Introduction

The construction of unbiased estimators and estimation of biases of estimators for several functions of unknown theoretical distribution is a great problem in statistics, especially in the case, when the given sample is not very large.

This means that the solution of the following general problems is needed:

$1^{\circ}$ Let $X$ be a sample of size n, and $t(X)$ a statistic. Then it is required to find the expression

$$E\,t(X) = \sum_{i=0}^{\infty} A_i\, n^{-i}. \tag{1}$$

Here $A_0$ is the leading term of the expectation, and the term $A_i$ defines the bias of order i, i=1,... .

$2^{\circ}$ Let $\tau$ be a non-random function of theoretical distribution (the parameter). Let $X$ be a sample of size n. It is required to construct a family of estimations $B_i(X)$, i=1,2,..., fulfilling the following conditions

$$\tau - EB_i(X) = O(n^{-i}), \quad i=1,2,\dots. \tag{2}$$

Then $B_i(X)$ is a $n^{-i}$ - biased estimation for $\tau$ .

We shall consider the special case when $t$ and $\tau$ are entire rational functions of the sample (correspondingly theoretical) moments of the arbitrary finite order.

Using the fact that for most functions, used in theoretical investigations and practical purposes, there exist approximations in the class of entire rational functions of moments(received, for instance, with the help of Taylor-type expansions of these functions), the result is rather general.

It is a well known fact that for a function $\tau(\cdot)$ of theoretical moments the simplest estimation is the corresponding function $t(\cdot)$ of the sample moments, but, in general, the estimation $t(\cdot)$ has a considerable bias, especially in the case when the order of the moments is large ( $>2$).

There exists a lot of solutions of problems $1^{\circ}$ and $2^{\circ}$

for special functions of moments, but in most cases some important restrictions are supposed, for instance:

1) The parent distribution is known.
2) The parent distribution belongs to a standard class of distributions (for example, to class $N(m, \delta)$).
3) The function $\tau$ has a special form.
4) The order of moments is not high ( $\leqslant 4$).

But there does not exist any general methodology for solving Problems $1^o$ and $2^o$ that might be realized in the form of standard packages and could be usable for reseachers.

The aim of the present paper is to solve Problems $1^o$ and $2^o$ for the class of entire rational functions of moments. All the solutions are analytical, and the formulae have the form, effectively realizable with the help of computer.

The methodological foundation for all the constructions will be based on the concept of the partition of integer, see [2], the necessary results will be given in Paragraph 2. Some useful preliminary results are given in [1].

## 2. The partition of an integer

Let $h$ be an integer. Then the vector $\quad \upsilon = (\upsilon_1, \dots, \upsilon_\Lambda)$, fulfilling the condition

$$\sum_{i=1}^{\Lambda} \upsilon_i = h \ , \quad \upsilon_i \in N \ , \ \upsilon_i > 0, \tag{3}$$

is said to be a partition of integer $h$; the components $\upsilon_i$ are parts of $\upsilon$, and $\Lambda = \Lambda(\upsilon)$ is the number of parts of $\upsilon$.

For quarranting the uniqueness of the representation the parts are assumed to be ordered:

$$\upsilon_i \geqslant \upsilon_{i+1} \ , \quad i = 1, \dots, \Lambda - 1 \ . \tag{4}$$

Notice that the partition is defined by its set of parts $\{ \upsilon_1, \dots, \upsilon_\Lambda \}$ . For every vector $w$ , having the components equal to the parts $\upsilon_1, \dots, \upsilon_\Lambda$ (in arbitrary order) there exists sugh a permutation P that $Pw = \upsilon$ . Sometimes we shall use some partitions $w$ , not fulfilling condition (4), as well.

Let $\upsilon^h$ denote the set of all different partitions $\upsilon^j = ( \upsilon_1^j, \dots, \upsilon_{\Lambda_j}^j )$ of $h$ , $j = 1, \dots, p(h); \ p(h)$ - the number of different partitions of $h$ is said to be the partition function of $h$.

Let $\upsilon_1 \in U^h, \upsilon_2 \in U^t$. Then, evidently, for $w = (\upsilon_1 : \upsilon_2) = = (\upsilon_1^1, \dots, \upsilon_{\Lambda_1}^1, \upsilon_1^2, \dots, \upsilon_{\Lambda_2}^2)$ holds: $w \in U^{h+t}$. Similarily, for a $\kappa \in N$ and $\upsilon \in U^h$ we have:

4

$$W = (\underbrace{v_1,\ldots,v_1}_{\kappa}, \ldots, \underbrace{v_\Lambda,\ldots,v_\Lambda}_{\kappa}), \quad w \in U^{\kappa h}.$$

If $h$ is fixed, we shall use the notation $U = U^h$.

Let $v^1, v^2 \in U$ . If $\Lambda_1 > \Lambda_2$ , then it is sometimes possible to get the partition $v^2$ by adding together some parts of $v^1$ . For every pair of partitions $v^1, v^2$ the coefficient $\beta(v^1, v^2)$ equals to the <u>number of different possibilities of getting</u> $v^2$ from $v^1$ .

<u>Example 1</u> Let $h = 5$, $v^1 = (2,1,1,1)$, $v^2 = (3,2)$. Then there are 4 different possibilities to get $v^2$ from $v^1$:

1) $v_1^2 = v_1^1 + v_2^1$ ; $v_2^2 = v_3^1 + v_4^1$ ;
2) $v_1^2 = v_1^1 + v_3^1$ ; $v_2^2 = v_2^1 + v_4^1$ ;
3) $v_1^2 = v_1^1 + v_4^1$ ; $v_2^2 = v_2^1 + v_3^1$ ;
4) $v_1^2 = v_2^1 + v_3^1 + v_4^1$ ; $v_2^2 = v_1^1$ ;

that means, $\beta(v^1, v^2) = 4$. Evidently, $\beta(v^2, v^1) = 0$.

For every $v \in U$ the set $U(v)$ is defined in the following way: $U(v) = \{v' : v' \in U, \beta(v, v') \neq 0\}$ and for every integer $\Lambda$ , $\Lambda \leqslant h$ the set $U_\Lambda$ is following: $U_\Lambda = \{v : v \in U, \Lambda(v) = \Lambda\}$. Let $U_\Lambda(v) = U_\Lambda \cap U(v)$ . From the definitions the following simple corollary is deduced.

<u>Corollary 1</u>

Let $U$ be the set of partitions of integer $h$ , $\Lambda \leqslant h$ , $v \in U$. Then the following equations hold:

$$U_\Lambda(v) = \{v\} ,$$
$$\beta(v, v) = 1 .$$

### 3. The <u>univariate moments and their products</u>

Let $X$ be a given random variable with distribution P. Assume that all theoretical moments, needed in our constructions, exist. Let $X$ be a sample of $X$ with size n, $X = (x_1, \ldots, x_n)$, $x_i \sim P$, independent. We shall use following notations:

$$\mu_p = E\, x^p , \tag{5}$$

$$m_p = \frac{1}{n} \sum_{i=1}^{n} x_i^p . \tag{5'}$$

Let $v$ be a partition, $v = (v_1, \ldots, v_\Lambda)$ . Then it defines uniquely a product of moments,

$$\mu(v) = \prod_{i=1}^{\Lambda} \mu_{v_i} , \tag{6}$$

$$m(v) = \prod_{i=1}^{\Lambda} m_{v_i} = n^{-\Lambda} \sum_{i_1=1}^{n} x_{i_1}^{v_1} \ldots \sum_{i_\Lambda=1}^{n} x_{i_\Lambda}^{v_\Lambda} . \tag{6'}$$

For every fixed distribution P, sample size n, given integers h and k we define the following classes of homogeneous polynomials of moments:

2

$$\mathcal{M}(h,k) = \left\{ \sum_{v \in \mho^h} \mu(v) \sum_{j=0}^{k} n^{-j} c(v,j) , \quad c(v,j) \in R \right\}, \qquad (7)$$

$$\mathfrak{m}(h,k) = \left\{ \sum_{v \in \mho^h} m(v) \sum_{j=0}^{k} n^{-j} c(v,j) , \quad c(v,j) \in R \right\}, \qquad (7')$$

where the coefficients $c(v,j)$ are independent from the distribution P and sample size n.

Let us denote $\mathcal{M}(h, \infty) = \mathcal{M}(h)$, $\mathfrak{m}(h, \infty) = \mathfrak{m}(h)$.
Evidently for $k_1 < k_2$ we have

$$\tau \in \mathcal{M}(h, k_1) \Rightarrow \tau \in \mathcal{M}(h, k_2). \qquad (8)$$

Let $t_i \in \mathfrak{m}(h_i)$ and $\tau_i \in \mathcal{M}(h_i)$. Then for finite k $\tau = \sum_{i=1}^{k} \tau_i$ and $t = \sum_{i=1}^{k} t_i$ are the entire rational functions of theoretical and, correspondingly, sample moments.
Let us denote the sets by $\mathfrak{m}$ and $\mathcal{M}$:

$$\mathfrak{m} = \left\{ t : t = \sum_{i=1}^{k} t_i , \ t_i \in \mathfrak{m}(h_i), \ h_i, k \in \mathcal{N}, \ i = 1, \dots, k \right\}, \qquad (9)$$
$$\mathcal{M} = \left\{ \tau : \tau = \sum_{i=1}^{k} \tau_i , \ \tau_i \in \mathcal{M}(h_i), \ h_i, k \in \mathcal{N}, \ i = 1, \dots, k \right\}.$$

In the paper the problems $1°$ and $2°$ will be solved for the functions $\tau \in \mathcal{M}(h)$ and $t \in \mathfrak{m}(h)$. Using the equalities (9) and the linearity of expectation, all results, proved for the classes $\mathfrak{m}(h)$ or $\mathcal{M}(h)$, hold for the classes $\mathfrak{m}$, $\mathcal{M}$, otherwise.

### 4. The expectation of the product of sample moments

**Theorem 1.** Let $m(v)$ be a product $(6')$. Then its expectation can be expressed in the following way:

$$E\,m(v) = \sum_{v' \in \mho(v)} \beta(v,v') n^{\delta - \delta'} F(\Lambda') \mu(v') , \qquad (10)$$

where $\Lambda = \Lambda(v)$, $\Lambda' = \Lambda(v')$ and the function $F(x)$ $(x \in \mathcal{N})$ is defined in the following way:

$$F(x) = \begin{cases} 1, & \text{if } x = 1, \\ (1 - \frac{1}{n}) \dots (1 - (x-1)/n), & \text{if } x > 1, \end{cases}$$

where $F(x)$ can be expressed in the following way:

$$F(x) = \sum_{i=0}^{x-1} f_i^x n^{-i} .$$

Proof The product $(6')$ consists of $n^{\delta}$ terms, having the form

$$n^{-\delta} x_{i_1}^{v_1} \dots x_{i_\Lambda}^{v_\Lambda} \qquad (i_j = 1, \dots, n; \ j = 1, \dots, \delta). \qquad (11)$$

The value of the expectation of the term (11) depends on the number of equations $i_j = i_g$ $(j, g = 1, \dots, \delta, j \neq g)$ between the indices. If $i_j = i_g$, then instead of product $x_{i_j}^{v_j} x_{i_g}^{v_g}$ the term $x_{i_j}^{v_j + v_g}$ may be written. That means, we have the expectation of (11) equal to $n^{-\delta} \mu(v')$, where $v' \in \mho(v)$.

6

The number of terms, having such expectation, equals to
$B(v, v')n(n-1)\ldots(n-(\wedge'-1)) = B(v, v') n^{-\wedge'} F(\wedge')$. Substitution
of all terms (11) by their expressions gives the formula (10).

Simple transformation of the formula (10) gives the expression of coefficients in the form (1):

Corollary 2

The expectation of the product of sample moments (6')
has the following expression:
$$E\, m(v) = \sum_{j=0}^{\wedge-1} n^{-j} \sum_{u=0}^{j} \mathcal{f}_{j-u}^{\wedge-u} \sum_{v' \in \mathcal{V}_{\wedge-u}(v)} B(v,v') \mu(v'). \tag{11}$$
Whereas the terms $\mathcal{f}_i^{x}$ and $B(v, v')$ are independent
from the distribution P and sample size $n$, we have proven
the result
$$E\, m(v) \in \mathcal{M}(h, \wedge-1). \tag{11'}$$

## 5. The expectation of entire rational function of sample moments

Using the connection (9) and the linearity of the expectation from (11') follows the accuracy of the formula (1)
for arbitrary function $t(m) \in \mathcal{M}(h)$. Hence the formula (1)
holds for arbitrary function $t \in \mathcal{M}$.

The problem $1^{\circ}$ is solved.

For the practical purposes it is useful to find the
exact values of some terms of the expectations. Most of them
may be simply deduced from the formula (11); we shall present
them in the form of corollaries.

Corollary 3

If $h$ is fixed and $t \in \mathcal{M}(h)$, then $E\,t \in \mathcal{M}(h)$.

From Corollaries 1 and 2 follows the well-known result:

Corollary 4
$$E\, m(v) = \mu(v) + \mathcal{O}(n^{-1}). \tag{12}$$

The formula (12) holds evidently for a more general case:

Corollary 5

Let $t$ be the entire rational function, $t = t(m)$, where
$m$ is a vector of sample moments. Then
$$E\,t = t(\mu) + \mathcal{O}(n^{-1}), \tag{13}$$
where $t(\mu)$ is the same function of responding theoretical
moments.

For practice it is useful to find the expression of the
bias of the first order of the estimation $m(v)$ for $\mu(v)$.
From formula (11) follows

## Corollary 6

$$E m(v) - \mu(v) = n^{-1} \left[ \sum_{v' \in U_{h-1}(v)} \beta(v,v') \mu(v') - 0.5 (\lambda^2 - \lambda) \mu(v) \right] + O(n^{-2}).$$

## 6. Estimation of sample central moments

One of the most frequently used functions from $m(h)$ is the sample central moment $\omega_h$ ,

$$\omega_h = \frac{1}{n} \sum_{j=1}^{n} \left( x_j - \frac{1}{n} \sum_{i=1}^{n} x_i \right)^h. \qquad (14)$$

After transformation of $\omega_h$ ,

$$\omega_h = \frac{1}{n} \sum_{j=1}^{n} \sum_{g=0}^{h} C_h^g (-1)^g n^{-g} x_j^{h-g} \left( \sum_{i=1}^{n} x_i \right)^g, \qquad (15)$$

we see that all terms in expression (15) respond to partitions of a special type, $v(g) = (h-g, 1, \ldots, 1)$ $(g = 0, 1, \ldots, h-1)$. Substituting the expectations $E m(v(g))$, calculated by formula (9), into expression (15) and ordering the result by powers of $n$ , we get the following

## Corollary 7

The expectation of the sample central moment $\omega_h$ is following:

$$E(\omega_h) = \sum_{j=0}^{h-1} n^{-j} \sum_{g=1}^{h-1} \sum_{u=0}^{j} \sum_{v' \in U_{g+1-u}(v(g))} (-1)^g \hat{C}_h^g \beta(v(g),v') + \frac{g+1-u}{j-1} \mu(v'), \qquad (16)$$

where the coefficient $\hat{C}_h^g$ is defined in the following way

$$\hat{C}_h^g = \begin{cases} 1, & \text{if } g = 0, \\ h!/(g!(h-g)!), & \text{if } 1 \le g \le h-2, \\ h-1, & \text{if } g = h-1. \end{cases}$$

From the Theorem 2 it is simple to get the leading term of the bias of the estimator $\omega_h$ for corresponding theoretical central moment $v_h = E(x - Ex)^h$ :

## Corollary 8

$$E \omega_h - v_h = n^{-1} \left\{ \sum_{g=1}^{h-2} (-1)^g C_h^g g \left[ \mu(v(g-1)) - 0.5 (g+1) \mu(v(g)) \right] + \right.$$

$$+ \sum_{g=2}^{h-2} (-1)^g C_h^g \cdot 0.5 (g^2 - g) \mu(v(g,2)) + (-1)^{h-1} 0.5 (h-1)^2 h \left[ \mu(v(h-2)) - \right.$$

$$\left. - \mu(v(h-1)) \right] \right\} + O(n^{-2}),$$

where $v(g,2) = (h-g, 2, 1, \ldots, 1)$.

Analogically, the higher moments and central moments of the statistic $\omega_h$ can be calculated; it is evident that

$$(\omega_h)^l \in m(hl), \quad E(\omega_h)^l \in \mathcal{M}(hl).$$

## 7. Estimation of the product $\mu(v)$ with bias of given order k

After the problem $1^0$ has been theoretically solved in paragraph 5, and most useful practical conclusions deduced

in paragraph 6, we shall consider the solution of problem $2^{\circ}$, proving the following

Theorem 2. There exists an estimator $M_k$ for the parameter $\tau \in \mathcal{Ul}(h)$ and given integer k, fulfilling the condition

$$E M_k - \tau = O(n^{-k-1}), \quad \ell \in N. \tag{18}$$

Proof bases on the idea of mathematical induction.

Let us assume that there exists an estimator $M_{i-1}$, fulfilling (18) for $k = i-1$, $M_{i-1} \in \mathcal{M}(h)$. Then by Corollary 3, $E M_{i-1} \in \mathcal{Ul}(h)$ and hence the difference $E M_{i-1} - \tau$ belongs to the class $\mathcal{Ul}(h)$ too; hence it has the following form

$$E M_{i-1} - \tau = \sum_{j=\ell}^{\infty} n^{-j} \sum_{v \in v} g(j,v) \mu(v).$$

Let us suppose that in the expression only $h-1$ first terms differ from zero, that means we have

$$E M_{i-1} = \tau + \sum_{j=\ell}^{i+h-2} n^{-j} \sum_{v \in v} g(j,v) \mu(v). \tag{19}$$

The next estimator $M_i$ will be constructed with the help of the first term of the expression (19), where the theoretical moments will be replaced by their sample analogues:

$$M_i = M_{i-1} - n^{-i} \sum_{v \in v} g(i,v) m(v). \tag{20}$$

For calculating the expectation of $M_i$ we use the formulae (11) and (12):

$$E m(v) = \mu(v) + \sum_{\ell=1}^{h-1} n^{-\ell} \sum_{v' \in v} \mu(v') c^v(\ell,v),$$

and get

$$E M_i = \tau + \sum_{j=i}^{i+h-2} n^{-j} \sum_{v \in v} g(j,v) \mu(v) - n^{-i} \sum_{v \in v} g(i,v) \mu(v) -$$
$$- n^{-i} \sum_{v \in v} g(i,v) \sum_{\ell=1}^{h-1} n^{-\ell} \sum_{v' \in v} \mu(v') c^v(\ell,v'),$$

that means,

$$E M_i = \tau + n^{-i-1} \left\{ \sum_{j=0}^{h-2} n^{-j} \sum_{v \in v} \mu(v) \left[ g(j+i+1,v) - \sum_{v \in v} g(i,v') c^{v'}(j+1,v) \right] \right\}, \tag{21}$$

where $g(i; h-1, v) = 0$ and if $\Lambda(v) < h$ then $c^v(j,v) = 0$ for $j \geqslant \Lambda(v)$.

The formula (21) demonstrates that in our assumption the connection (18) holds for $\ell = i$ as well, whereas the expression has the form (19).

For completing the proof notice that from Corollary 4 it follows that the equation (18) holds in the case k=0; then for the estimator $M_c$ may be chosen the correspondent function of sample moments $M_c = \tau(m)$. Theorem 2 is proven.

9

## Corollary 10

For $\mu(v)$ the estimator

$$M_1(v) = m(v)\left(1 + (\lambda^2 - \lambda)/(2n)\right) - \frac{1}{n}\sum_{v' \in \mathcal{V}_{h-1}(v)} \beta(v,v')\, m(v')$$

is $n^{-2}$-biased,

$$E M_1(v) - \mu(v) = \mathcal{O}(n^{-2}).$$

Similarily with the help of Corollary 7 the $n^{-2}$-biased estimation of the k-th central moment $\mathcal{V}_k$ may be constructed . For this purpose in the equation (17) the theoretical moments must be replaced by corresponding sample analogs.

## 8. The existence of the unbiased estimator $M_\infty$

The following question is the existence of an unbiased estimation $M_\infty = \lim_{k \to \infty} M_k$ , having the property

$$E M_\infty = \lim_{k \to \infty} E M_k = \tau. \tag{22}$$

Let us use the notation

$$E M_i - \tau = \sum_{j=i+1}^{i+h-1} n^{-j} \sum_{v \in \mathcal{V}} \mu(v)\, g^i(j,v).$$

For convergence (22) the following series must be convergent:

$$\sum_{i=0}^{\infty} \sum_{j=i+1}^{i+h-1} n^{-j} \sum_{v \in \mathcal{V}} \mu(v)\, g^i(j,v). \tag{23}$$

In fact, the series (23) is a linear combination of $p(h)(h-1)$ different series (some of them may be constantly equal to zero), having the following form:

$$\sum_{i=0}^{\infty} n^{-i} g^i(j+i,v), \qquad j = 1,\ldots,h-1;\ v \in \mathcal{V}. \tag{24}$$

For the absolute convergence of series (23) it is necessary and sufficient that all series (24) $(j = 1,\ldots,h-1;\ v \in \mathcal{V})$ are absolutely convergent.

For analysing the convergence of series (24) we regard the expression of coefficient $g^i(j,v)$ by previous coefficients $g^+(j',v')$ $(+ \le i,\ j' = 1,\ldots,h-1,\ v' \in \mathcal{V})$. From (21) we get:

$$g^i(j+i,v) = g^{i-1}(j+i+1,v) - \sum_{v' \in \mathcal{V}} g^{i-1}(i,v')\, c^{v'}(j+1,v), \tag{25}$$

this means that the coefficient $g^i(j,v)$ is a linear transformation of previous coefficients, where the coefficients $c^{v'}(j+1,v)$ of the transformation are independent from the index i.

For investigating the convergence of the series (24)

10

let us construct the majorant series

$$g = \sum_{i=0}^{\infty} n^{-i} g_i ,$$ 

(26)

where $\quad g_i = \max_{\substack{1 \le j \le h-1 \\ v \in V}} |q^i(j+i,v)| ,$

from the convergence of the series $g$ the convergence of all series (24), and hence the convergence of series (23) follows.

From (25) we have, using $[1]$ ,

$$g_i \le g_{i-1} (\sum_{v \in V} c^v(j+1,v)+1) \le g_{i-1} (B(h)+1) ,$$

where $B(h)$ is the Bell's number, responding to number $h$ , see $[2]$ .

Let us suppose that the sample size $n$ fulfills the condition

$$n > B(h)+1 .$$

(27)

Then we have the inequality

$$(g_i n^{-i})/(g_{i-1} n^{-i+1}) = (B(h)+1)/n = c < 1 ,$$

hence the series (25) and (23) are convergent.

In fact we have proved the following

<u>Theorem 3</u>.For the existence of unbiased estimator $M_\infty$ for parameter $\tau \in \mathcal{M}(h)$ it is sufficient that the sample size $n$ fulfills the condition (27). The estimation $M_\infty$ may be constructed as the sum of convergent series (23) with the typical coefficient (25).

Notice that the existence of estimation depends only on the number $h$ (premised all $EX^j, j \le h$ exist),but not on the values of moments or the concrete form of distribution P.

If $\tau = \mu(v)$ , then instead of $h$ the value $\Delta(v)$ may be used, as it follows from (21).

Whereas all constants $c^{v'}(j,v)$ may be simply calculated from (11), the exact calculation of estimator $M_\infty$ is only a technical problem.

9. <u>Estimation of variance of a function</u> $t \in \mathcal{M}(h)$

Let $t \in \mathcal{M}(h)$ , then $t^2 \in \mathcal{M}(2h)$ . Suppose the given distribution is of order $2h$ , then by Corollary 3 $Et^2 \in \mathcal{M}(2h)$.

Let us regard the term $(Et)^2$ :

$$(Et)^2 = \sum_{v \in V} \sum_{v' \in V} \mu(v) \mu(v') \left[ \sum_{j=0}^{h-1} n^{-j} g(v,j) \sum_{i=0}^{h-1} n^{-i} g(v',i) \right] =$$

$$= \sum_{w \in \vartheta^{2k}} \mu(w) \sum_{j=c}^{2\Delta-2} n^{-d} q(w,j), \tag{28}$$

where [1])

$$g(w,j) = \begin{cases} \sum_{i=0}^{j} g(v,i)g(v',j-i), & \text{if } w=(v \vdots v'), \ v,v' \in \vartheta^h, \\ 0, & \text{otherwise.} \end{cases}$$

From (28) follows the inclusion $(Et)^2 \in \mathcal{M}(2h)$, consequently $\mathcal{D}t \in \mathcal{M}(2h)$, and from Theorems 2 and 3 and Corollary 8 follows

<u>Corollary 11.</u> Let $t \in \mathcal{M}(h)$. Then $\mathcal{D}t \in \mathcal{M}(2h)$ and for every k there exists statistic $T_k \in \mathcal{M}(2h)$, being $n^{-k}$-biased estimation of $\mathcal{D}t$.

<u>Corollary 12.</u> For sample size n,

$$n > \beta(2h)+1, \tag{29}$$

there exists an unbiased estimator $T_\infty$ for $\mathcal{D}t$,

$$T_\infty = \lim_{k \to \infty} T_k,$$

where the estimators $T_k$ may be constructed step-wise, using the formula(19).

Let us calculate the leading term of variance of estimation $M_i$ for $\mu(v)$, $i=0,1,\dots$.

<u>Theorem 4</u> The variance of estimator $M_k$ is following

$$\mathcal{D}M_k = \prod_{i=1}^{\Delta} \mu_{v_i}^2 \left\{ \sum_{i=1}^{\Delta} \frac{\mu_{2v_i}}{\mu_{v_i}^2} - \Delta^2 \right\} \bar{n}^{-1} + O(n^{-2}), \quad (k=0,1,\dots). \tag{30}$$

<u>Proof</u> At first we shall prove the formula (30) for k=0, $M_0 = m(v)$. Let us denote $W = (v_1, v_1, \dots, v_\Delta, v_\Delta)$, where $v=(v_1,\dots,v_\Delta)$, then $\Delta(w)=2\Delta$. For calculation $\mathcal{D}m(v)$ we need to calculate the terms $E(m(v))^2$ and $(Em(v))^2$, using the fact

$$(m(v))^2 = m(w), \quad (\mu(v))^2 = \mu(w).$$

followingly,

$$\begin{cases} E(m(v))^2 = \mu(w) + n^{-1}\left\{ \sum_{w' \in \vartheta_{2\Delta-1}(w)} \beta(w,w')\mu(w') - (2\Delta^2-\Delta)\mu(w) \right\} + O(n^{-2}) \\ (Em(v))^2 = \mu(w) + n^{-1}\left\{ 2\mu(v) \sum_{v' \in \vartheta_{\Delta-1}(v)} \beta(v,v')\mu(v') - (\Delta^2-\Delta)(\mu(v))^2 \right\} + O(n^{-2}) \end{cases} \tag{31}$$

Without restriction of generality we may suppose all parts of partitions to be different, then $\beta(v,v')=1$ for all $v' \in \vartheta_{\Delta-1}(v)$ (real results may be received as a result of adding similar terms).

Then we divide the set $\vartheta_{2\Delta-1}(w)$ into two subsets:

1) In the article the partitions need not fulfil the condition(4).

12

$$\mathcal{V}_{\varkappa} = \left\{ (v_i + v_i, v_1, v_1, \ldots, v_{i-1}, v_{i-1}, v_{i+1}, v_{i+1}, \ldots, v_\Lambda, v_\Lambda) ; \ i = 1, \ldots, \Lambda \right\}; \quad \mathcal{X}(\mathcal{V}_{\varkappa}) = \Lambda,$$

$$\mathcal{V}_{\varkappa\varkappa} = \left\{ (v_i + v_j, v_1, v_1, \ldots, v_{i-1}, v_{i-1}, v_i, v_{i+1}, v_{i+1}, \ldots, v_{j-1}, v_{j-1}, v_j, v_{j+1}, v_{j+1}, \ldots, v_\Lambda, v_\Lambda) ; \right.$$
$$\left. i = 1, \ldots, \Lambda-1 ; \ j = i+1, \ldots, \Lambda \right\} : \ \mathcal{X}(\mathcal{V}_{\varkappa\varkappa}) = 0.5(\Lambda^2 - \Lambda).$$

The coefficients $B(w, w')$ are evidently following:

$$B(w, w') = \begin{cases} 1, & \text{if } w' \in \mathcal{V}_{\varkappa}, \\ 2, & \text{if } w' \in \mathcal{V}_{\varkappa\varkappa}. \end{cases}$$

From the definition we get the expression of set $\mathcal{V}_{\varkappa\varkappa}$ by $\mathcal{V}_{\Lambda-1}(v)$:

$$\mathcal{V}_{\varkappa\varkappa} = \left\{ (v : v'), v' \in \mathcal{V}_{\Lambda-1}(v) \right\}. \tag{32}$$

Using the expressions (31) and (32) we have

$$\mathcal{D} m(v) = n^{-1} \left\{ \sum_{w' \in \mathcal{V}_{2\Lambda-1}(w)} B(w, w') \mu(w') - 2 \sum_{v' \in \mathcal{V}_{\Lambda-1}(v)} \mu(v) \mu(v') - \Lambda^2 \mu(w) \right\} + O(n^{-2}) =$$

$$= n^{-1} \left\{ \sum_{i=1}^{\Lambda} \mu_{2v_i} \prod_{\substack{j=1 \\ j \neq i}}^{\Lambda} \hat{\mu}_{v_j}^2 - \Lambda^2 \prod_{i=1}^{\Lambda} \hat{\mu}_{v_j}^2 \right\} + O(n^{-2}). \tag{33}$$

The formula (32) is equivalent to the (29).

Let us regard now the estimator $M_1$.

$$M_1 = m(v) - \frac{1}{n} \left[ \sum_{v' \in \mathcal{V}(v)} B(v, v') m(v') - 0.5(\Lambda^2 - \Lambda) m(v) \right] + O(n^{-2}),$$

$$M_1^2 = m(w) - \frac{1}{n} \left[ \sum_{v' \in \mathcal{V}(v)} 2B(v, v') m(v) m(v') - (\Lambda^2 - \Lambda) m(w) \right] + O(n^{-2}),$$

$$E M_1^2 = E m(w) - \frac{1}{n} \left[ \sum_{v' \in \mathcal{V}(v)} 2B(v, v') \mu(v) \mu(v') - (\Lambda^2 - \Lambda) \mu(w) \right] + O(n^{-2}),$$

using the first expression from the formula (31) we get

$$E M_1^2 = \mu(w) + \frac{1}{n} \left[ \sum_{w \in \mathcal{V}^{\varkappa}} \mu(w) - \Lambda^2 \mu(w) \right] + O(n^{-2}), \tag{34}$$

but from the definition of $\bar{E} M_1^2$ it follows that

$$(E M_1)^2 = \mu(w) + O(n^{-2}). \tag{35}$$

Substituting the coefficient by $n^{-1}$ from (34) into the expression of $DM_1$, we receive the leading term of the variance of $M_1$ equal to that of $M_o$.

So as all the following estimations differ from $M_1$ only by higher terms, the result holds for all k.

Theorem 4 is proven.

Analogically to Theorem 4 it is possible to prove the following

Corollary 13

Let $t_i \in m(h_i, k_i), \ i = 1, \ldots, \ell$, $t = \prod_{i=1}^{\ell} t_i$.

4                    13

Then $t \in m(h,k)$, $h = \prod_{i=1}^{A} h_i$, $k = \prod_{i=1}^{A} k_i$ and if given distribution P is of order $h$, then $Et \in \mathcal{M}(h,k)$.

If $t$ is a sample analog tö a theoretical parameter $\tau$, then for the existence of the unbiased estimation $T_\infty$ for $\tau$ the fulfilling of the condition

$$n > B(h) + 1$$

is sufficient; then $T_\infty$ may be constructed as a result of the step-wise process, described in Theorem 2.

From Corollary 13 it follows that the results described in article [9] may be generalized for calculating arbitrary moments (central moments) of arbitrary entire rational function $t$ of sample moments, including unbiased estimation $T_\infty$ of given parameter $\tau$ (for sample size $n$, large enough for its existence).

Example 1. We study the products of moments, responding to all partitions of numbers 1,...,5. For the case of uniform distribution $\mathcal{U}(0,2)$ for all these products the exact values and expected biases for sample sizes 5,10 and 20 are calculated. The modelling experiment consists of trials of 10 000 (for $n = 10$ and 20) and 20 000 (for $n = 5$) samples; for every sample the all products of moments are computed and their average (over trials) is found. The difference between the average and exact value gives the empirical bias. The results of modelling experiment are given in Table 1.

Example 2. Here we study the product moment $\tau = \mu_1^3$ and calculate for it a series of $n^i$-biased estimations $M_i$ for $\tau$ ($i = 0,1,2,...$).

$M_0 = m_1^3$,

$EM_0 = \mu_1^3 + \frac{1}{n}(3\mu_2\mu_1 - 3\mu_1^3) + \frac{1}{n^2}(\mu_3 - 3\mu_2\mu_1 + 2\mu_1^3)$,

$M_1 = M_0 - \frac{1}{n}(3m_2m_1 - 3m_1^3)$,

$EM_1 = \mu_1^3 + \frac{1}{n^2}(-2\mu_3 + 9\mu_2\mu_1 - 7\mu_1^3) + \frac{1}{n^3}(3\mu_3 - 9\mu_2\mu_1 + 6\mu_1^3)$,

$M_2 = M_1 - \frac{1}{n^2}(-2m_3 + 9m_2m_1 - 7m_1^3)$,

Having calculated the expectation of the estimator $M_{i-1}$ with help of its term of order $n^{-i}$ the following estimator $M_i$ is defined, as it follows:

$M_3 = M_2 - \frac{1}{n^3}(-6m_3 + 21m_2m_1 - 15m_1^3)$,

$M_4 = M_3 - \frac{1}{n^4}(-14m_3 + 45m_2m_1 - 31m_1^3)$,

$M_5 = M_4 - \frac{1}{n^5}(-30m_3 + 93m_2m_1 - 63m_1^3)$,

etc.

14

Table 1

Modelling experiment for estimation of product of sample moments.

| N | Product of moments | N of terms | Exact value | Theoretical bias | | | Empirical bias | | | Bias of empirical bias | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | n=5 | n=10 | n=20 | n=5 | n=10 | n=20 | n=5 | n=10 | n=20 |
| 1 | $m_1$ | 1 | 1 | 0 | 0 | 0 | −0.0035 | −0.0022 | 0.0012 | −0.0035 | −0.0022 | 0.0012 |
| 2 | $m_2$ | 1 | 1.3333 | 0 | 0 | 0 | −0.0055 | −0.0039 | 0.0026 | −0.0055 | −0.0039 | 0.0026 |
| 3 | $m_1^2$ | 2 | 1 | 0.0667 | 0.0333 | 0.0167 | 0.0596 | 0.0297 | 0.0189 | −0.0071 | −0.0037 | 0.0022 |
| 4 | $m_3$ | 1 | 2 | 0 | 0 | 0 | −0.0087 | −0.0061 | 0.0048 | −0.0087 | −0.0061 | 0.0048 |
| 5 | $m_2 m_1$ | 2 | 1.3333 | 0.1333 | 0.0667 | 0.0333 | 0.1228 | 0.0613 | 0.0372 | −0.0105 | −0.0054 | 0.0039 |
| 6 | $m_1^3$ | 3 | 1 | 0.2 | 0.1 | 0.05 | 0.1891 | 0.0953 | 0.0532 | −0.0109 | −0.0047 | 0.0032 |
| 7 | $m_4$ | 1 | 3.2 | 0 | 0 | 0 | −0.0149 | −0.0097 | 0.0085 | −0.0149 | −0.0097 | 0.0085 |
| 8 | $m_3 m_1$ | 2 | 2 | 0.24 | 0.12 | 0.06 | 0.2235 | 0.1122 | 0.0664 | −0.0165 | −0.0078 | 0.0064 |
| 9 | $m_2^2$ | 2 | 1.7778 | 0.2844 | 0.1422 | 0.0711 | 0.2687 | 0.1352 | 0.0770 | −0.0157 | −0.0070 | 0.0059 |
| 10 | $m_2 m_1^2$ | 3 | 1.3333 | 0.3591 | 0.1787 | 0.0891 | 0.3431 | 0.1724 | 0.0940 | −0.0160 | −0.0063 | 0.0049 |
| 11 | $m_1^4$ | 4 | 1 | 0.4123 | 0.2032 | 0.1008 | 0.3968 | 0.1977 | 0.1048 | −0.0155 | −0.0055 | 0.0040 |
| 12 | $m_5$ | 1 | 5.3333 | 0 | 0 | 0 | −0.0268 | −0.0160 | 0.0148 | −0.0268 | −0.0160 | 0.0148 |
| 13 | $m_4 m_1$ | 2 | 3.2 | 0.4267 | 0.2133 | 0.1067 | 0.3994 | 0.2013 | 0.1173 | −0.0273 | −0.0120 | 0.0106 |
| 14 | $m_3 m_2$ | 2 | 2.6667 | 0.5333 | 0.2667 | 0.1333 | 0.5085 | 0.2569 | 0.1425 | −0.0248 | −0.0098 | 0.0092 |
| 15 | $m_3 m_1^2$ | 3 | 2 | 0.6240 | 0.3093 | 0.1540 | 0.5990 | 0.3004 | 0.1617 | −0.0250 | −0.0091 | 0.0077 |
| 16 | $m_2^2 m_1$ | 3 | 1.7778 | 0.6542 | 0.3236 | 0.1609 | 0.6305 | 0.3157 | 0.1679 | −0.0237 | −0.0069 | 0.0070 |
| 17 | $m_2 m_1^3$ | 4 | 1.3333 | 0.7019 | 0.3424 | 0.1690 | 0.6792 | 0.3354 | 0.1747 | −0.0227 | −0.0070 | 0.0057 |
| 18 | $m_1^5$ | 5 | 1 | 0.7280 | 0.3493 | 0.1708 | 0.7066 | 0.3433 | 0.1754 | −0.0214 | −0.0060 | 0.0048 |

For calculation of the unbiased estimation we must find the following sums (coefficients by $m_i^3$, $m_2 m_1$ and $m_3$):

$$a = 1 + \frac{3}{n} + \frac{7}{n^2} + \frac{15}{n^3} + \ldots \; ,$$

$$b = -\frac{3}{n} - \frac{9}{n^2} - \frac{21}{n^3} - \frac{45}{n^4} - \ldots \; ,$$

$$c = \frac{2}{n^2} + \frac{6}{n^3} + \frac{14}{n^4} + \frac{30}{n^5} + \ldots \; .$$

Notice, that the numerators of terms of series a are the Stirling's numbers $S(2,n)$ of the second order and series b and c are simple functions of series a:

$$a = 1 + \sum_{i=1}^{n} \frac{2^{i+1}-1}{n^i} = 1 + 2\frac{2}{n-2} - \frac{1}{n-1} = \frac{n^2}{(n-2)(n-1)} \; ,$$

$$b = -\frac{3}{n} a = \frac{-3n}{(n-2)(n-1)} \; ; \quad c = \frac{2}{n^2} a = \frac{2}{(n-2)(n-1)} \; ,$$

and, followingly,

$$M_\infty = (n^2 m_1^3 - 3n\, m_2 m_1 + 2 m_3)/((n-2)(n-1)).$$

It is simple to check that $EM_\infty = \mu_1^3$.

For existence of the estimation (convergence the series a, b and c) it is sufficient when $n > 3$.

The modelling experiment consists of series of 10 000 samples of size $n = 5$ and $n = 10$ from the uniform distribution $\mathcal{U}(0,1)$. For $k = 0, 1, \ldots, 10$ and $\infty$ the $n^{-k}$-biased estimations, their expected values and averages of series, also their variances by the series are calculated. The results are given in Table 2 (as all results became constant for $k = 8, \ldots, \infty$ they are not given in table).

Example 3. Let us estimate the variance of estimations $M_0$ and $M_\infty$ for $\mathcal{T} = \mu_1^3$ (from Example 3):
$$M_0 = m_1^3; \; M_0^2 = m_1^6; \; EM_0^2 = \mu_1^6 + n^{-1}\{15 \mu_2 \mu_1^4 - 15 \mu_1^6\} + O(n^{-2}),$$
$$EM_0 = \mu_1^3 + \frac{1}{n}(3\mu_2 \mu_1 - 3\mu_1^3) + O(n^{-2}),$$
$$(EM_0)^2 = \mu_1^6 + n^{-1}(6\mu_2 \mu_1^4 - 6\mu_1^6) + O(n^{-2}),$$
hence $\mathcal{D} M_0 = n^{-1}(9\mu_2 \mu_1^4 - 9\mu_1^6) + O(n^{-2})$;
$$M_\infty^2 = (n^4 m_1^6 - 6 n^3 m_2 m_1^4 + 9 n^2 m_2^2 m_1^2 + 4 n^2 m_3 m_1^3 - 12 n m_3 m_2 m_1 + 4 m_3^2)/((n-1)^2 (n-2)^2),$$
$$EM_\infty^2 = \frac{1}{(n-1)^2(n-2)^2}\{n^4[\mu_1^6 + n^{-1}(15\mu_2 \mu_1^4 - 15\mu_1^6) + O(n^{-2})] - 6n^3[\mu_2 \mu_1^4 + O(n^{-1})] + O(n^{-2})\}.$$

Whereas $EM_\infty = \mu_1^3$, then $(EM_\infty)^2 = \mu_1^6$, and we have
$$\mathcal{D} M_\infty = \frac{n^4}{(n-1)^2(n-2)^2}\{\mu_1^6 + n^{-1}(9\mu_2 \mu_1^4 - 15\mu_1^6) + O(n^{-2})\} - \mu_1^6 \; ,$$

16

Table 2

The $n^{-k}$-biased and unbiased estimations of $\mu_1^3$ in the case of $\mathcal{U}(0,1)$. The true value is $\mu_1^3 = 0,125$.

| $k$ | $n = 5$ | | | | $n = 10$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $M_k$ | $DM_k$ | theor. bias | bias of emp. bias | $M_k$ | $DM_k$ | theor. bias | bias of emp. bias |
| 0 | 0.15343 | 0.01168 | 0.02500 | 0.00343 | 0.137702 | 0.005345 | 0.0125 | 0.000202 |
| 1 | 0.130453 | 0.01125 | 0.00500 | 0.000453 | 0.126429 | 0.005142 | 0.00125 | 0.000179 |
| 2 | 0.126470 | 0.0110 | 0.00100 | 0.000470 | 0.125300 | 0.005105 | 0.000125 | 0.000175 |
| 3 | 0.125671 | 0.01104 | 0.00020 | 0.000471 | 0.125186 | 0.005098 | 0.000013 | 0.000173 |
| 4 | 0.125511 | 0.01101 | 0.00004 | 0.000471 | 0.125175 | 0.005097 | 0.000001 | 0.000174 |
| 5 | 0.125478 | 0.01100 | 0.000008 | 0.000470 | 0.125174 | 0.005096 | 0 | 0.000174 |
| 6 | 0.125472 | 0.01099 | 0.000002 | 0.000470 | 0.125174 | 0.005096 | 0 | 0.000174 |
| 7 | 0.12547 | 0.01099 | 0 | 0.000470 | 0.125174 | 0.005096 | 0 | 0.000174 |
| $\infty$ | 0.12547 | 0.01099 | 0 | 0.000470 | 0.125174 | 0.005096 | 0 | 0.000174 |

so as $\dfrac{n^4}{(n-1)^2(n-2)^2}\mu_1^6 = \mu_1^6 + \dfrac{6}{n}\mu_1^6 + O(n^{-2})$,

we have $\partial M_\infty = n^{-1}(9\mu_2\mu_1^4 - 9\mu_1^6) + O(n^{-2})$.

In the case of uniform distribution $\mathcal{U}(0,1)$ the leading term of variances $DM_0$ and $DM_\infty$ equals to 0.0469 for $n = 5$ and 0.0094 for $n = 10$. The result is in good agreement with the computational experiment (see Table 2), where the sequence of empirical variances $DM_k$ of estimators is almost constant.

References

I. Тийт Э.-М. Вычисление математических ожиданий некоторых функций эмпирических моментов. Труды ВЦ ТГУ, 1986, 53, 60-85.
2. Эндрюс Г. Теория разбиений. М., 1982.

### НЕСМЕЩЕННЫЕ И $n^k$-СМЕЩЕННЫЕ ОЦЕНКИ ЦЕЛО-РАЦИОНАЛЬНЫХ ФУНКЦИЙ МОМЕНТОВ
#### Э.Тийт

Р е з ю м е

Для цело-рациональной функции эмпирических моментов $t$ строится оценка математического ожидания в виде разложения по степеням объема выборки $n$ (см. формулу (I)). В таблице I приведены результаты соответствующего моделирующего эксперимента.

Пусть $\tau$ цело-рациональная функция теоретических моментов, притем предполагается, что все необходимые теоретические моменты существуют. Для $\tau$ строится оценка $M_k$ со смещением порядка $n^{-k-1}$, к=I,2, ... (см. формулу (I8)) и доказывается, что при достаточно большом объеме выборки $n$ (см. формулу (27)) существует несмещенная оценка $M_\infty$, выражаемая в виде сходящего ряда (22). Доказывается, что дисперсии оценок $M_k$ (к=0, I, ...) и $M_\infty$ имеют равные главные члены (см. формулу (30)). В таблице 2 приведены результаты моделирующего эксперимента, иллюстрирующие полученные результаты.

# ON THE STABILITY OF k-MEANS CLUSTERING
## IN METRIC SPACES

### K. Pärna

Summary

Let P be a probability measure on a separable metric space $(T,d)$. Let $\mathcal{A}^*_k(P)$ be the class of all minimizing sets for the clustering criterion $W(A,P)$ (see formula (1.1)), and let $W_k(P)=\inf\{W(A,P):|A|=k\}$ (each $A^*=\{a_1^*,\ldots,a_k^*\}$ from the class $\mathcal{A}^*_k(P)$ will be called the optimal k-centre for the measure P). Let $\{P_n\}$ be a sequence of probability measures, weakly converging to P, $P_n \Rightarrow P$. In Part I the convergence $W_k(P_n) \longrightarrow W_k(P)$ and some other results are proven, under certain restrictions on $\varphi$, P, and $\{P_n\}$. These statements serve as generalizations of our previous results (see [2], Theorem). In Part II the case of separable Hilbert spaces is considered. Conditions are found to ensure the convergence $A_n^* \to \mathcal{A}_2^*(P)$ (in certain sense), for any sequence $\{A_n^*\}$, $A_n^* \in \mathcal{A}_2^*(P_n)$.

### Part I

## 1. Introduction

Consider a separable metric space $(T,d)$ together with a probability measure P on it. Define the clustering criterion by

$$W(A,P) = \int \min_{1 \le i \le k} \varphi(d(x,a_i))P(dx),{}^{*)} \quad A = \{a_1,\ldots,a_k\} \subset T, \quad (1.1)$$

where the function $\varphi$ satisfies the following restrictions:

1) $\varphi:[0,\infty) \longrightarrow [0,\infty)$,
2) $\varphi$ is continuous,
3) $\varphi$ is nondecreasing,
4) $\varphi(r) = 0 \Longleftrightarrow r=0$,
5) there exists a constant $\lambda$, such that $\varphi(2r) \le \lambda \varphi(r)$, $r \ge 0$ ( $\Delta_2$-property).

These restrictions remain in force during the paper. Further, let

$$W_k(P)=\inf\{W(A,P):|A|=k\}, \quad (1.2)$$
$$\mathcal{A}^*_k(P) = \{A:W(A,P)=W_k(P), |A|=k\}. \quad (1.3)$$

Definition 1. Any $A^* = \{a_1^*,\ldots, a_k^*\}$ from the class $\mathcal{A}^*_k(P)$

---

*) The domain of integration is T everywhere, if not specified.

is said to be an optimal k-centre for the measure P.

Any optimal k-centre provides the best approximation of P by discrete distribution concentrated at k points. Also, each $A^*$ generates an optimal partition of T to k exhaustive and mutually exclusive regions (clusters) $S_1^*,\ldots,S_k^*$, satisfying

$$S_i^* \subseteq \left\{ x: \psi(d(x,a_i^*)) \leq \psi(d(x,a_j^*)), \text{ for each } j, j \neq i \right\}$$

- hence the term 'k-means clustering'.

In our earlier paper [2] the problem of consistency of the criterion's infimum $W_k(P)$ was considered. To be more precise, let $x_1,\ldots,x_n$ be a random sample of n independent observations on the distribution P. The measure $P_n$, obtained by placing mass $1/n$ at each of $x_1,\ldots,x_n$, is called the empirical measure. In the paper noted above the consistency of $W_k(P_n)$, i.e., the almost sure convergence

$$W_k(P_n) \longrightarrow W_k(P) \tag{1.4}$$

has been proven, under some restrictions on P. Earlier Pollard has shown the same for the space $R^m$ [3].

In this part of the paper the convergence (1.4) will be shown to hold for arbitrary sequence $\{P_n\}$, provided that $P_n \Rightarrow P$ (weakly). For the case of $T = R^m$ an analogous result has been given by Abaya and Wise [1].

It is well known that the topology of the weak convergence is metrizable by Prokhorov metrics $\pi$ [6]. Thereby, our result can be reformulated in stability terms: small changes in P (in the sense of $\pi$) do not have a significant effect on the infimum value $W_k(P)$. It follows that 'good' k-centres for some measure P' are 'good' for the measure P too, provided that $\pi(P,P')$ is small.

## 2. The stability theorem for the infimum of $W(A,P)$

Let us introduce some restrictions on the probability measures P and $P_n$, $=1,2,\ldots$, considered as being given on the separable metric space T:

C1) $P_n \Rightarrow P$,

C2) for some $y_o \in T$ the function $\psi(d(x,y_o))$ is uniformly integrable with respect to $\{P_n\}$,

C3) the strong inequalities

$$W_1(P) > W_2(P) > \ldots > W_k(P) \tag{2.1}$$

hold.

The following theorem is one of the main assertions of this work.

Theorem 1. Under the restrictions C1)-C3) the convergence

$$W_k(P_n) \longrightarrow W_k(P), \quad n \longrightarrow \infty,$$

takes place.

Proof. The proof of this theorem does not make a significant difference from the proof of the related theorem in our previous work [2] where the particular case of empirical measures $P_n$ was considered. Both proofs make use of the convergence

$$W(A, P_n) \longrightarrow W(A, P) \qquad (2.2)$$

with certain fixed A, $|A| = k$, and the only difference is how to ground this convergence. In the case of empirical measures (2.2) follows directly from the strong law of large numbers (SLLN). In the general case of $\{P_n\}$ satisfying C1) and C2) this convergence is a result of applying the theory of weak convergence (Theorem 5.4 from [6] can be exploited). To avoid tedious repetitions, we shall not present the details of the proof here. Merely, we give some lemmas, serving as main steps of that proof. These lemmas will be used in the further sections.

Lemma 1. Under the conditions C1)-C3) there exists a sphere $B(x_0, M)$ which contains all the sets from the classes $\mathcal{A}_k^*(P_n)$, $n \geqslant 1$.

A similar result holds for the class

$$\mathcal{A}_k^\varepsilon(P) = \{A: W(A, P) < W_k(P) + \varepsilon, \ |A| = k\}$$

- the class of the '$\varepsilon$-optimal k-centres for the measure P'.

Lemma 2. Let us propose that $W_{k-1}(P) > W_k(P)$ and $\int \varphi(d(x, y_0)) P(dx) < \infty$ for some $y_0 \in T$. Then, for each $\varepsilon$, $0 < \varepsilon < W_{k-1}(P) - W_k(P)$, there exists a ball $B(x_0, M)$ which contains all the sets from the class $\mathcal{A}_k^\varepsilon(P)$.

To the spheres $B(x_0, M)$ in these two lemmas the following Lemma 3 applies.

Lemma 3. Under the conditions C1 and C2), for any fixed ball $B(z, R) \subset T$, the uniform convergence

$$\sup_{A \subset B(z, R)} |W(A, P_n) - W(A, P)| \to 0, \quad n \longrightarrow \infty,$$

takes place.

This lemma can be proven by means of a uniform convergence theorem given by Ranga Rao (see Theorem 3.2 in [4]). In a special case when $P_n$ is the empirical measure Lemma 3 coincides with Lemma 1 from [2].

From the results given above one useful conclusion can

21

6

be drawn.

Corollary 1. Let the conditions C1) - C3) be satisfied and let $A_n^*$ be an optimal k-centre for the measure $P_n$. Then
$$W(A_n^*, P) \longrightarrow W_k(P), \quad n \longrightarrow \infty .$$

Proof. Due to Lemma 1 all $A_n^*$ are contained in some sphere $B(x_0, M)$. So we have
$$\left| W(A_n^*, P) - W_k(P) \right| \leqslant$$
$$\leqslant \left| W(A_n^*, P) - W(A_n^*, P_n) \right| + \left| W(A_n^*, P_n) - W_k(P) \right| \leqslant$$
$$\leqslant \sup_{A \subset B(x_0 M)} \left| W(A, P) - W(A, P_n) \right| + \left| W_k(P_n) - W_k(P) \right| .$$

Now it remains to apply Lemma 3 and Theorem 1 to obtain the needed result.

### 3. The case of empirical measures.

In this section it will be shown that the restrictions C1) and C2) are weak enough to include the important case of empirical measures $P_n$, corresponding to the measure $P$. In such a way it becomes clear that Theorem 1 is a direct generalization of Theorem in [2].

Lemma 4. Let $\int \varphi(d(x, y_0)) \, P(dx) < \infty$ for some $y_0 \in T$. Then the empirical measures $\{P_n\}$ satisfy (with probability 1) the conditions C1) and C2).

Proof. It has been proven by Varadarajan [5] that if $P_n$ is the empirical measure corresponding to $P$, then $P_n \Longrightarrow P$ almost surely (a.s.). So C1) is satisfied with probability 1. To show this for C2) let
$$I_{\alpha, n} = \int_{\varphi(d(x, y_0)) \geqslant \alpha} \varphi(d(x, y_0)) \, P_n(dx) \tag{3.1}$$

and let $I_\alpha$ be the same for $P$, $I_\alpha < \infty$. We have to prove the uniform integrability of $\varphi(d(x, y_0))$, i.e., the relation
$$\lim_{\alpha \to \infty} \sup_n I_{\alpha, n} = 0, \quad \text{a.s.} \tag{3.2}$$

Suppose (controversially) that
$$\lim_{\alpha \to \infty} \sup_n I_{\alpha, n} = M > 0. \tag{3.3}$$

Then two cases are possible: $M = \infty$ or $M < \infty$. If $M = \infty$ then (3.3) implies $\sup_n I_{\alpha, n} = \infty$, for each $\alpha > 0$.

On the other hand $I_{\alpha, n} \to I_\alpha < \infty$ (a.s.), according to SLLN. Hence $M = \infty$ with probability 0.

Let $M < \infty$. Due to $I_\alpha \downarrow 0 \ (\alpha \to \infty)$ there exists a num-

22

ber $\alpha_0$ such that

$$I_\alpha < \frac{M}{2} \ , \quad \alpha > \alpha_0. \tag{3.4}$$

On the other hand, from (3.3) it follows that, for any $\alpha \geqslant 0$, $\sup_n I_{\alpha,n} \geqslant M$, and thus, for any $\alpha \geqslant 0$, there exists an integer $n_\alpha$ such that

$$I_{\alpha,n_\alpha} \geqslant \frac{M}{2} \ . \tag{3.5}$$

Let $\alpha_i = \alpha_0 + i$, $i = 1,2,\ldots$ and let $n_i$ be the least index with property

$$I_{\alpha_i,n_i} \geqslant \frac{M}{2} \ , \quad i=1,2,\ldots \tag{3.6}$$

Then, in view of simple inequalities

$$I_{\alpha_i,n_{i+1}} \geqslant I_{\alpha_{i+1},n_{i+1}} \geqslant \frac{M}{2} \ ,$$

it is clear that $n_{i+1} \geqslant n_i$ $(i=1,2,\ldots)$. Furthermore, the sequence $\{n_i\}$ can not be bounded. Indeed, if $n_i = N$, for every $i$ greater that some $i_0$, then

$$\lim_{i\to\infty} I_{\alpha_i,n_i} = \lim_{i\to\infty} I_{\alpha_i,N} = 0,$$

which contradicts (3.6). Hence $n_i \longrightarrow \infty$ , as $i \longrightarrow \infty$. Now, from the inequalities

$$I_{\alpha_1,n_i} \geqslant I_{\alpha_i,n_i} \geqslant \frac{M}{2} \ , \quad i \geqslant 1 ,$$

it follows that

$$\lim_i I_{\alpha_1,n_i} \geqslant \frac{M}{2} \ .$$

At the same time, by the SLLN, almost surely

$$\lim_i I_{\alpha_1,n_i} = I_{\alpha_1} ,$$

which is less than $\frac{M}{2}$ (see formula (3.4)) - a contradiction. Hence the case $M < \infty$ also occurs with probability 0.

This proves the lemma.

## 4. Optimal and $\mathcal{E}$-optimal partitions.

In this section the reader's attention is turned to the partitions of T rather than k-centres.

Let $\mathcal{P}_k(T)$ be the class of all measurable k-partitions of the metric space T,

$$\mathcal{P}_k(T) = \Big\{ S:S=\{S_1,\ldots,S_k\}, \ \emptyset \neq S_i \in \mathcal{T}, \ \bigcup_{i=1}^{k} S_i = T, \ S_i \cap S_j = \emptyset, \ i \neq j \Big\}.$$

Now, let us introduce the functional

$$W(A,S,P) = \sum_{i=1}^{K} \min_{a_j \in A} \int_{S_i} \varphi(d(x,a_j)) P(dx), \ A=\{a_1,\ldots,a_k\}, S \in \mathcal{P}_k(T) \tag{4.1}$$

23

- the summary measure of the goodness of the approximation of the region $S_i$ by means of the best point from the given set A.

Definition 2. Any point $a(s_i) \in T$, which minimizes the integral (over a)

$$W(a,S_i,P) = \int_{S_i} \varphi(d(x,a)) \, P(dx),$$

is called the centre of the region $S_i$.

Definition 3. Any set of the form $A(S) = \{a(S_1),\ldots, a(S_k)\}$ with $a(S_i)$ being a centre of $S_i$, is called the k-centre of the partition $S = \{S_1,\ldots,S_k\}$. The class of all k-centres for the partition S will be denoted by cent S.

It follows directly from the Definitions 2 and 3 that any $A(S)$ is optimal in the sense of $W(\cdot,S,P)$, i.e., for any $A(|A|=k)$ and $S, S \in \mathcal{G}_k(T)$, the inequality

$$W(A,S,P) \geqslant W(A(S),S,P) \, , \, A(S) \in \text{cent } S, \tag{4.2}$$

holds. Now, putting

$$W(S,P) = \inf \{W(A,S,P) : |A|=k\}, \tag{4.3}$$

(4.2) implies

$$W(S,P) = W(A(S),S,P), \tag{4.4}$$

provided that cent $S \neq \emptyset$. According to (4.1) the latter formula reduces to

$$W(S,P) = \sum_{i=1}^{k} \inf_{a} \int_{S_i} \varphi(d(x,a)) \, P(dx). \tag{4.5}$$

Definition 4. Let $A = \{a_1,\ldots, a_k\} \subset T$. Any partition of the form

$$S^{\varphi}(A) = \{S_1^{\varphi}(A),\ldots, S_k^{\varphi}(A)\},$$

where

$$S_i^{\varphi}(A) \subseteq \{x : \varphi(d(x,a_i)) \leq \varphi(d(x,a_j)) \text{ for all } j, j \neq i\}, \text{ is}$$

said to be a generalized Dirichlet partition with respect to A.

If $\varphi$ is strictly increasing, the latter definition gives us the minimum distance partition (also known as Dirichlet or Voronoi partition).

From Definition 4 it is clearly seen that any partition $S^{\varphi}(A)$ is optimal in the sense of $W(A,\cdot,P)$, i.e.,

$$W(A,S^{\varphi}(A),P) = \inf \{W(A,S,P) : S \in \mathcal{G}_k(T)\}. \tag{4.6}$$

Also, the formula (1.1) may be rewritten as

$$W(A,P) = \sum_{i=1}^{k} \int_{S_i^{\varphi}(A)} \varphi(d(x,a_i)) P(dx) \equiv W(A, S^{\varphi}(A), P). \tag{4.7}$$

Finally, from the relations (1.2) and (4.3) – (4.7) it fol-

lows that
$$W_k(P)=\inf_A W(A,P)=\inf_S W(S,P)= \inf_{A,S} W(A,S,P) \qquad (4.8)$$
where $|A|=k$, $S \in \mathcal{Y}_k(T)$.

Definition 5. Say a k-partition S is optimal with respect to the measure P if
$$W(S,P)=W_k(P),$$
and $\mathcal{E}$-optimal if
$$W(S,P) < W_k(P) +\mathcal{E}.$$
It is seen from (4.7) and (4.3) that
$$W(A,P)=W(A,S^{\mathcal{Y}}(A),P) \geqslant \inf_{|A|=k} W(A,S^{\mathcal{Y}}(A),P) \equiv W(S^{\mathcal{Y}}(A),P),$$
and, thereby, the $\mathcal{E}$-optimality of A implies the optimality of each $S^{\mathcal{Y}}(A)$.

Now we will give the main result of this section. The following theorem states that 'good' partitions can not have the regions with arbitrary small P-measure.

Theorem 2. Let $W_{k-1}(P) > W_k(P)$ and $\int \psi(d(x,y_0))P(dx) < \infty$ for some $y_0 \in T$. Then for each $\mathcal{E}$, $0 < \mathcal{E} < W_{k-1}(P) - W_k(P)$, there exists an $\alpha > 0$ such that the inequalities
$$P(S_i) \geqslant \alpha, \quad i=1,\ldots,k,$$
hold for all $\mathcal{E}$-optimal k-partitions $S=\{S_1,\ldots, S_k\}$.

Proof. Suppose (controversially) that for some $\mathcal{E}_0$, $0 < \mathcal{E}_0 < W_{k-1}(P)-W_k(P)$, there exist the $\mathcal{E}_0$-optimal partitions $S^{(n)}= \{S_1^{(n)},\ldots, S_k^{(n)}\}$, $n=1,2,\ldots$, such that for some region, say $S_k^{(n)}$, the inequality
$$P(S_k^{(n)}) < \frac{1}{n}$$
holds. The main idea now is to show that combining $S_k^{(n)}$ with any other region of $S^{(n)}$ we get a (k-1)-partition without any significant enlarging of $W(S,P)$.

Define the (k-1)-partition
$$\widetilde{S}^{(n)}= \{S_1^{(n)},\ldots, S_{k-2}^{(n)}, S_{k-1}^{(n)} \cup S_k^{(n)}\}.$$
Due to the formula (4.5)
$$W(\widetilde{S}^{(n)},P)-W(S^{(n)},P)=\inf_a \int_{S_{k-1}^{(n)} \cup S_k^{(n)}} \psi(d(x,a))\, P(dx)-$$
$$-\inf_a \int_{S_{k-1}^{(n)}} \psi(d(x,a))P(dx)-\inf_a \int_{S_k^{(n)}} \psi(d(x,a))P(dx) . \qquad (4.9)$$

For the possible nonexistence of the optimal centre for the region $S_{k-1}^{(n)}$ let us introduce an $\mathcal{J}$-optimal centre $a_{\mathcal{J}}(S_{k-1}^{(n)})$,

7

satisfying

$$\int_{S_{k-1}^{(n)}} \psi(d(x,a_\gamma(S_{k-1}^{(n)})))P(dx) < \inf_a \int_{S_{k-1}^{(n)}} \psi(d(x,a))P(dx)+\delta , \quad (4.10)$$

with

$$\delta = \frac{\Delta}{2km},$$

$$\Delta = W_{k-1}(P)- W_k(P) - \mathcal{E}_o > 0 .$$

It is not too hard to see that the relations (4.9) and (4.10) imply

$$W(\tilde{S}^{(n)},P)-W(S^{(n)},P)< \int_{S_k^{(n)}} \psi(d(x,a_\gamma(S_{k-1}^{(n)})))P(dx)+\delta . \quad (4.11)$$

We now show that there exists a sphere $B(x_o,M)$ which contains all $a_\gamma(S_{k-1}^{(n)})$, $n=1,2,\ldots$ Let $a_\gamma(S_i^{(n)})$ be a $\delta$-optimal centre for $S_i^{(n)}, i=1,\ldots,k$, and let $cent_\gamma(S^{(n)}) = \{a_\gamma(S_i^{(n)})\}_{i=1}^k$. For the inequalities

$$W(cent_\gamma S^{(n)},S^{(n)},P) < W(S^{(n)},P) + k\cdot\delta \le$$

$$\le W_k(P)+ \mathcal{E}_o + k\delta = W_k(P) + \mathcal{E}_o +\frac{\Delta}{2n}$$

and

$$W(cent_\gamma S^{(n)},P) \le W(cent_\gamma S^{(n)}, S^{(n)},P)$$

(the latter follows from (4.6) and (4.7)) the set $cent_\gamma S^{(n)}$ is an $(\mathcal{E}_o+\frac{\Delta}{2n})$-optimal k-centre for the measure P. Since for any $n=1,2,\ldots$

$$\mathcal{E}_o+\frac{\Delta}{2n} < \mathcal{E}_o+\Delta < W_{k-1}(P) - W_k(P),$$

Lemma 2 applies and hence the sphere $B(x_o,M)$ exists which contains all the sets $cent_\delta S^{(n)}$, $n \ge 1$, including the points $a_\gamma(S_{k-1}^{(n)})$.

Let us return to (4.11) now. For the triangle inequality we have

$$\int_{S_k(n)} \psi(d(x,a_\gamma(S_{k-1}^{(n)})))P(dx)+\delta \le \int_{S_k^{(n)}} \psi(d(x,x_o)+M)P(dx)+\delta . \quad (4.12)$$

As the last integral is finite (see Appendix, Lemma A1), the property of the absolute continuity of the Lebesgue integral implies that there exists a $\gamma>0$ such that

$$\int_{S_k^{(n)}} \psi(d(x,x_o)+M)P(dx) < \Delta - \frac{\Delta}{2k} \le \Delta - \delta , \quad (4.13)$$

for any region $S_k^{(n)}$ satisfying

$$P(S_k^{(n)}) < \gamma. \quad (4.14)$$

Now, let n be an integer greater than $1/\gamma$. Then, by the presumption that $P(S_k^{(n)}) < 1/n$, the inequality (4.14) is fulfilled, and the relations (4.11)-(4.13) readily give us

$$W(\widetilde{S}^{(n)}, P) - W(S^{(n)}, P) < \Delta \quad . \qquad (4.15)$$

Consequently,

$$W(S^{(n)}, P) > \; W(\widetilde{S}^{(n)}, P) - \Delta \geqslant W_{k-1}(P) - \Delta = W_k(P) + \mathcal{E}_o \; ,$$

i.e., the partition $S^{(n)}$ is not $\mathcal{E}_o$-optimal for P. We reached the contradiction.

This completes the proof of Theorem 2.


## Part II


Our aim here is to obtain some convergence results for $\{A_n^*\}$ - the sequence of optimal k-centres for $P_n$ - provided that $P_n \Rightarrow P$. Only the case of separable Hilbert spaces will be considered.

### 5. Some preliminary results.

It is well known that all real separable Hilbert spaces are isometrically isomorphic to the real space $\ell_2$. So we may take $T = \ell_2$ at once, and let $X = (X_1, X_2, \ldots)$ be a random element with values in $\ell_2$. Assume that X has the distribution P of the second order, i.e.

$$\sum_{j=1}^{\infty} EX_j^2 = E \| X \|^2 = \int \| x \|^2 P(dx) < \infty \; . \qquad (5.1)$$

This part of the paper deals with a particular case of the clustering criterion $W(A, P)$, namely

$$W(A,P) = \int \min_{1 \leqslant i \leqslant k} \| x - a_i \|^2 P(dx), \quad A = \{a_1, \ldots, a_k\} \subset \ell_2. \qquad (5.2)$$

Further, let $S = \{S_1, \ldots, S_k\}$ be a Dirichlet partition of $\ell_2$, generated by the set $A = \{a_1, \ldots, a_k\}$ :

$$S_i \subseteq \{x : \| x - a_i \| \leqslant \| x - a_j \| \text{ for all } j, \; j \neq i\} \; .$$

Then it is not difficult to check the relations

$$W(A,P) = \sum_{i=1}^{k} \int_{S_i} \| x - a_i \|^2 P(dx) =$$

$$= E \| X \|^2 + \sum_{i=1}^{k} P(S_i) \cdot ( \| a(S_i) - a_i \|^2 - \| a(S_i) \|^2), \qquad (5.3)$$

where $a(S_i) = E_{S_i}(X)$ - the conditional mean of X in the region $S_i$. In other words, $a(S_i)$ is the centre of $S_i$.

Also, for arbitrary partition $S = \{S_1, \ldots, S_k\}$ of $\ell_2$ the formulae

**7***

$$W(S,P) = \sum_{i=1}^{k} \int_{S_i} \| x - a(S_i) \|^2 \, P(dx),$$

$$W(S,P) = E\|X\|^2 - \sum_{i=1}^{k} P(S_i) \, \| a(S_i) \|^2 \qquad (5.4)$$

are valid.

Now, let us introduce the subspaces

$$\mathcal{L}_n = \left\{ x : x = (x_1, \ldots, x_n, 0, 0, \ldots) \right\} \subset \ell_2,$$

and $\quad \mathcal{L}_{(n)} = \left\{ x : x = (0, \ldots, 0, x_{n-1}, x_{n+2}, \ldots) \right\} \subset \ell_2,$

with corresponding projectors $\mathbb{T}_n$ and $\mathbb{T}_{(n)}$:

$$\mathbb{T}_n x = (x_1, \ldots, x_n, 0, \ldots) \in \mathcal{L}_n,$$

$$\mathbb{T}_{(n)} x = (0, \ldots, 0, x_{n+1}, x_{n+2}, \ldots) \in \mathcal{L}_{(n)}, \quad x \in \ell_2.$$

Let $P_n$ and $P_{(n)}$ be the distributions of (random) projections $\mathbb{T}_n X$ and $\mathbb{T}_{(n)} X$, correspondingly, that is

$$P_n = P \mathbb{T}_n^{-1}, \qquad (5.5)$$

$$P_{(n)} = P \mathbb{T}_{(n)}^{-1}. \qquad (5.6)$$

In the following lemma the notation $\delta_j^2$ for the variance of $X_j$ is used.

Lemma 6. Let P be a probability measure of the second order on $\ell_2$. Then, for any $A = \{ a_1, \ldots, a_k \} \subset \ell_2$, and for any $n = 1, 2, \ldots$, the inequalities

$$W(A,P) \geq W(\mathbb{T}_n A, P_n) + W(\mathbb{T}_{(n)} A, P_{(n)}), \qquad (5.7)$$

$$W_k(P) \leq W_k(P_n) + \sum_{j=n+1}^{\infty} \delta_j^2 \qquad (5.8)$$

hold.

Proof. Observing that

$$\| y \|^2 = \| \mathbb{T}_n y \|^2 + \| \mathbb{T}_{(n)} y \|^2, \quad y \in \ell_2,$$

from (5.2) it follows that

$$W(A,P) = \int \min_i \left[ \| \mathbb{T}_n x - \mathbb{T}_n a_i \|^2 + \| \mathbb{T}_{(n)} x - \mathbb{T}_{(n)} a_i \|^2 \right] P(dx). (5.9)$$

Since the minimum of the sum is never less than the sum of the minimums, the right side of (5.9) is no less than the sum of two integrals, which is equivalent to (5.7).

To prove (5.8), note first that it is sufficient to consider the case $EX = 0$, only. Indeed, any quantity in (5.8) does not depend on the mean of the distribution P. Thus, for any $A = \{ a_1, \ldots, a_k \}$, the simple relations

$$W(\widetilde{\Pi}_n A, P) \equiv \int \min_i \| x - \widetilde{\Pi}_n a_i \|^2 P(dx) =$$

$$= \int \min_i \| \widetilde{\Pi}_n x - \widetilde{\Pi}_n a_i \|^2 P(dx) + \int \sum_{j=n+1}^{\infty} x_j^2 P(dx) =$$

$$= W(\widetilde{\Pi}_n A, P_n) + \sum_{j=n+1}^{\infty} \delta_j^2 \tag{5.10}$$

hold. By taking into account that
$$W_k(P) \leqslant W(\widetilde{\Pi}_n A, P)$$
and
$$W_k(P_n) = \inf_{|A|=k} W(\widetilde{\Pi}_n A, P_n),$$

we obtain (5.8).

The proof is completed.

From the inequalities (5.7) and (5.8) it is seen that

$$W(A,P) - W_k(P) \geqslant W(\widetilde{\Pi}_{(n)} A, P_{(n)}) - \sum_{j=n+1}^{\infty} E X_j^2 . \tag{5.11}$$

## 6. The main lemma.

For convenience, let us denote
$$\| X \|^2_{(n)} = \| \widetilde{\Pi}_{(n)} x \|^2 = \sum_{j=n+1}^{\infty} x_j^2 , \quad x \in \ell_2.$$

**Lemma 7.** Let $P$ be a probability measure of the second order on $\ell_2$. Then, for any $\varepsilon > 0$, there exists a number $n$ such that any $A \in \mathcal{A}_k^\varepsilon (P)$ contains at least one point $a_i$ statisfying

$$\| a_i \|_{(n)} \leqslant 2\sqrt{2\varepsilon} . \tag{6.1}$$

**Proof.** First write that, for any $\varepsilon > 0$, any $A = \{ a_1, \ldots, a_k \}$ and every $n \geqslant 1$,

$$W(\widetilde{\Pi}_{(n)} A, P_{(n)}) = \int \min_i \| x - \widetilde{\Pi}_{(n)} a_i \|^2 P_{(n)}(dx) =$$

$$= \int \min_i \| \widetilde{\Pi}_{(n)} x - \widetilde{\Pi}_{(n)} a_i \|^2 P(dx) \geqslant$$

$$\geqslant 2\varepsilon \cdot P \{ \min_i \| \widetilde{\Pi}_{(n)} x - \widetilde{\Pi}_{(n)} a_i \|^2 \geqslant 2\varepsilon \} =$$

$$= 2\varepsilon \cdot P \{ \min_i \| x - a_i \|_{(n)} \geqslant \sqrt{2\varepsilon} \}. \tag{6.2}$$

Suppose, controversially, that for some $\varepsilon > 0$ there exists, for every $n$, an $\varepsilon$-optimal set $A = \{ a_1, \ldots, a_k \}$ satisfying

$$\| a_i \|_{(n)} > 2\sqrt{2\varepsilon} \quad \text{for every } i = 1, \ldots, k. \tag{6.3}$$

Observing that under the condition (6.3) the set
$\{ x : \min_i \| x - a_i \|_{(n)} \geqslant \sqrt{2\varepsilon} \}$ contains the other set

29

8

$\left\{x: \|x\|_{(n)} \leqslant \sqrt{2\mathcal{E}}\right\}$, (6.2) implies that

$$W(\widetilde{\mathbb{1}}_{(n)}A, P_{(n)}) \geqslant 2\mathcal{E} \cdot P\left\{\|x\|_{(n)} \leqslant \sqrt{2\mathcal{E}}\right\}.$$

Furthermore, by means of Chebyshev's inequality, we get

$$W(\widetilde{\mathbb{1}}_{(n)}A, P_{(n)}) \geqslant 2\mathcal{E} \cdot (1 - \frac{E\|X\|_{(n)}^2}{2\mathcal{E}}) = 2\mathcal{E} - E\|X\|_{(n)}^2 , \qquad (6.4)$$

where X is a random element with distribution P. But P is of the second order and thus, for some large N, we have

$$E\|X\|_{(n)}^2 = \sum_{j=n+1}^{\infty} EX_j^2 < \frac{\mathcal{E}}{2} , \text{ if } n > N. \qquad (6.5)$$

Let us fix an $n_0 > N$ and let $A_0$ be an $\mathcal{E}$-optimal set satisfying (6.3) (with $n=n_0$). Combining (6.5) with (6.4) we see that

$$W(\widetilde{\mathbb{1}}_{(n_0)}A_0, P_{(n_0)}) > \mathcal{E} + E\|X\|_{(n_0)}^2. \qquad (6.6)$$

Show now that this inequality contradicts to the $\mathcal{E}$-optimality of $A_0$. Indeed, putting $A=A_0$ and $n=n_0$ in (5.11) it is seen, by (6.6), that

$$W(A_0, P) > W_k(P) + \mathcal{E}.$$

This proves the Lemma.

## 7. The case k=2.

In this section it will be shown that in a special case when k=2, Lemma 7 can be strenghtened up to the assertion that all (both) points of $\mathcal{E}$-optimal 2-centres have a small 'tail'. In what follows, X is considered as a random element of $\ell_2$ with the distribution P.

Let $a(S_i) = (a_1(S_i), a_2(S_i), \ldots) \in \ell_2$ be the coordinate-wise presentation of the centre of the set $S_i \subset \ell_2$, $a(S_i) = E_{S_i}(X)$.

**Lemma 8.** Let $E\|X\|^2 < \infty$ and $EX=0$. Then, for any $\mathcal{E}$ satisfying $0 < \mathcal{E} < W_1(P) - W_2(P)$, there exists a number $\beta > 0$ such that the inequalities

$$\frac{1}{\beta}|a_j(S_2)| \leqslant |a_j(S_1)| \leqslant \beta|a_j(S_2)|, \quad j=1,2,\ldots \qquad (7.1)$$

hold simultaneously for all $\mathcal{E}$-optimal 2-partitions $S=\{S_1, S_2\}$ ($|\cdot|$ denotes the absolute value).

**Remark.** In this Lemma we assume that $W_1(P) > W_2(P)$. It is known, due to Lemma 2 in [2], that this inequality is satisfied whenever P is not concentrated at any single point.

**Proof.** Let $S=\{S_1, S_2\}$ be an $\mathcal{E}$-optimal partition of $\ell_2$ into 2 regions. Then, by the obvious fact that

$$a(S_1) \cdot P(S_1) + a(S_2) \cdot P(S_2) = EX = 0,$$

we have
$$a_j(S_1)P(S_1)+a_j(S_2)P(S_2)=0, \text{ for every } j=1,2,\ldots \qquad (7.2)$$
By Theorem 2 there exists $\alpha > 0$ such that $P(S_i) \geqslant \alpha$, $i=1,2$. Thus we may write
$$\left|a_j(S_1)\right| = \left|a_j(S_2)\right| \frac{P(S_2)}{P(S_1)} .$$
Since
$$\frac{\alpha}{1-\alpha} \leqslant \frac{P(S_2)}{P(S_1)} \leqslant \frac{1-\alpha}{\alpha}$$
it suffices to take $\beta = (1-\alpha)/\alpha$ to obtain (7.1).

**Lemma 9.** Let $E \|X\|^2 < \infty$ and $EX=0$. Then, for any $\mathcal{E}_0$ satisfying $0 < \mathcal{E}_0 < W_1(P) - W_2(P)$, there exist $\beta_c > 0$ and $C_0 > 0$ such that inequalities
$$\|a_1\|_{(n)} \leqslant \beta_0 \|a_2\|_{(n)} + C_0\sqrt{\mathcal{E}} , \qquad n=1,2,\ldots \qquad (7.3)$$
hold simultaneously for all $A = \{a_1, a_2\} \in \mathcal{A}_2^{\mathcal{E}}(P)$, provided that $\mathcal{E} \leqslant \mathcal{E}_0$ . These inequalities remain in force after changing the roles of $a_1$ and $a_2$.

**Proof.** Take $\mathcal{E} = \mathcal{E}_c$ in Theorem 2 and Lemma 8 and let $\beta_0 > 0$ and $\alpha_0 > 0$ be the corresponding values of $\beta$ and $\alpha$. Let $A = \{a_1, a_2\} \in \mathcal{A}_2^{\mathcal{E}}(P)$ and let $S = \{S_1, S_2\}$ be a Dirichlet partition generated by $A$. It is clear that $A \in \mathcal{A}_2^{\mathcal{E}_c}(P)$ whenever $\mathcal{E} \leqslant \mathcal{E}_0$ and, according to Lemma 8,
$$\|a(S_1)\|_{(n)} \leqslant \beta_0 \|a(S_2)\|_{(n)} , \qquad n=1,2,\ldots \qquad (7.4)$$
On the other hand, the formulae (5.3) and (5.4) directly imply that
$$W(A,P)-W(S,P) = \sum_{i=1}^{2} P(S_i) \|a(S_i) - a_i\|^2 . \qquad (7.5)$$
Since $W(A,P)-W(S,P) \leqslant W(A,P)-W_k(P) \leqslant \mathcal{E}$ (due to $\mathcal{E}$-optimality of $A$), we have
$$\sum_{i=1}^{2} P(S_i) \|a(S_i)-a_i\|^2 \leqslant \mathcal{E} .$$
Observing that $P(S_i) \geqslant \alpha_0$, $i=1,2$, one can obtain
$$\|a(S_i)-a_i\| \leqslant \sqrt{\frac{\mathcal{E}}{\alpha_0}} , \qquad i=1,2. \qquad (7.6)$$
By the simple inequalities
$$\left| \|a(S_i)\|_{(n)} - \|a_i\|_{(n)} \right| \leqslant \|a(S_i)-a_i\|_{(n)} \leqslant \|a(S_i)-a_i\| ,$$
(7.6) implies
$$\left| \|a(S_i)\|_{(n)} - \|a_i\|_{(n)} \right| \leqslant \sqrt{\frac{\mathcal{E}}{\alpha_c}} , \qquad i=1,2, \qquad n=1,2,\ldots$$
Combining this with (7.4) it is easily seen that
$$\|a_1\|_{(n)} \leqslant \beta_c \|a_2\|_{(n)} + \frac{1}{\alpha_c} \sqrt{\frac{\mathcal{E}}{\alpha_0}} .$$

**8***

Now the substitution $C_o = 1/\alpha_o\sqrt{\alpha_o}$ completes the proof.

The following Lemma is the main result of this section.

**Lemma 10.** Let $E\|X\|^2 < \infty$ , $EX=0$ and $0 < \mathcal{E}_o < W_1(P) - W_2(P)$. Then, for any $\mathcal{E}$ satisfying $0 < \mathcal{E} \leqslant \mathcal{E}_o$ , there exists a number n such that inequalities

$$\| a_i \|_{(n)} \leqslant C\sqrt{\mathcal{E}}, \quad i=1,2, \tag{7.7}$$

where C depends only on $\mathcal{E}_o$, hold simultaneously for all $A=\{a_1,a_2\} \in \mathcal{A}_2^{\mathcal{E}}(P)$.

<u>Proof.</u> Let us have an $\mathcal{E}_o$, $0 < \mathcal{E}_o < W_1(P)-W_2(P)$. By Lemma 7, for any $\mathcal{E}$, including the case when $\mathcal{E} \leqslant \mathcal{E}_o$ , there exists a number n such that any $A=\{a_1,a_2\}$ from $\mathcal{A}_2^{\mathcal{E}}(P)$ contains at least one point (say $a_2$) satisfying

$$\| a_2 \|_{(n)} \leqslant 2\sqrt{2\mathcal{E}}. \tag{7.8}$$

Thus (7.7) is shown to be valid for $a_2 \in A$. To show the same for $a_1$, we apply Lemma 9. From the inequalities (7.3) and (7.8) it is seen that

$$\| a_1 \|_{(n)} \leqslant \beta_o 2\sqrt{2\mathcal{E}} + C_o\sqrt{\mathcal{E}}.$$

Now it becomes clear that (7.7) holds for both, $a_2$ and $a_1$, with

$$C=\max\{2\sqrt{2}, \ 2\sqrt{2}\beta_o + C_o\}.$$

Thus the proof is completed.

## 8. Main results.

Now nearly everything has been done to formulate some theorems concerning the convergence of the $P_n$-optimal 2-centres to the P-optimal 2-centre in Hilbert spaces. However, the type of convergence is not specified as yet.

<u>Definition 6.</u> The quantity

$$h(A,B)=\max\left\{\sup_{a\in A}\inf_{b\in B} d(a,b), \ \sup_{b\in B}\inf_{a\in A} d(a,b)\right\}$$

is called the Hausdorf distance between the sets A and B.

<u>Definition 7.</u> The quantity

$$h(A,\mathcal{A})= \inf_{B\in\mathcal{A}} h(A,B)$$

is called the Hausdorf distance from the set A to the class of sets $\mathcal{A}$.

<u>Theorem 3.</u> Let a probability measure P of the second order in space $\ell_2$ be given. Let $A_n=\{a_1^n, a_2^n\}$, n=1,2,... be a minimizing sequence for $W(A,P)$, i.e., $W(A_n,P) \to W_2(P)$. Then

1) there exists a subsequence $\{A_{n'}\}$ such that, for some $A \in \mathcal{A}_2^*(P)$, $h(A_{n'},A) \to 0$;

2) $h(A_n, \mathcal{A}_2^*(P)) \to 0$, $n \to \infty$ . $\tag{8.1}$

32

Proof. Show first that 1) implies 2). Indeed, if 2) does not hold, then, for some $\mathcal{E} > 0$, there exists a subsequence $\{A_{n}\}$ such that

$$h(A_{n'}, \mathcal{A}_2^*(P)) > \mathcal{E}, \quad \text{for all } n'. \tag{8.2}$$

Since $\{A_{n'}\}$ is a minimizing sequence, it contains, by assertion 1), a further subsequence $\{A_{n''}\}$ which converges to some $A \in \mathcal{A}_2^*(P)$. Then, clearly,

$$h(A_{n''}, \mathcal{A}_2^*(P)) \longrightarrow 0, \quad n'' \longrightarrow \infty,$$

in contradiction to (8.2). Hence 1) implies 2).

To prove assertion 1), we may assume, as in Lemma 6, that the measure P is centered. Let us fix some arbitrary labelling of the elements in each $A_n = \{a_1^n, a_2^n\}$ ($n \geq 1$), and let $B_1 = \{a_1^1, a_1^2, \ldots\}$ – the set of the first elements of $A_n$. Our aim is to show that $B_1$ is relatively compact in $\ell_2$. Then $B_1$ will contain a converging (in norm) sequence.

It is well known ([7], page 52) that a set M from the complete separable metric space T (including the case $T = \ell_2$) is relatively compact if, for any $\delta > 0$, there exists a relatively compact $\delta$-net. Obviously, it suffices to regard the $\delta$-s small enough.

To demonstrate the relative compactness of $B_1$, let us fix $\mathcal{E}_o$, $0 < \mathcal{E}_o < W_1(P) - W_2(P)$. Then, by Lemma 10, there exists, for any $\mathcal{E} \leq \mathcal{E}_o$, a positive integer $m_1$ such that the inequalities

$$\| a_1^n \|_{(m_1)} \leq c \sqrt{\mathcal{E}} \quad , \quad i = 1, 2,$$

with C depending on $\mathcal{E}_o$ (and not on $\mathcal{E}$), hold for all $\mathcal{E}$-optimal $A_n = \{a_1^n, a_2^n\}$. Since $W(A_n, P) \rightarrow W_2(P)$, it is clear that the set $A_n$ is $\mathcal{E}$-optimal whenever n exeeds some $N = N(\mathcal{E})$. On the other hand, there is a positive integer $m_2$ such that

$$\| a_i^n \|_{(m_2)} \leq c \sqrt{\mathcal{E}} , \quad i = 1, 2,$$

for all n from 1 to N. Let $m = \max\{m_1, m_2\}$. Then the inequalities

$$\| a_i^n \|_{(m)} \leq c \sqrt{\mathcal{E}} , \quad i = 1, 2,$$

hold for all $\mathcal{E} \leq \mathcal{E}_o$, and all $n \geq 1$.

Now we shall indicate a relatively compact $\delta$-net for $B_1$, provided that $\delta \leq \delta_o = c \sqrt{\mathcal{E}_o}$. Since

$$\| a_1^n - \widetilde{\pi}_m a_1^n \| = \| a_1^n \|_{(m)} \leq c \sqrt{\mathcal{E}} , \quad n \geq 1,$$

then putting $\delta = c \sqrt{\mathcal{E}}$ the set $\widetilde{\pi}_m B_1 = \{\widetilde{\pi}_m a_1^1, \widetilde{\pi}_m a_1^2, \ldots\} \subset \mathcal{L}_m$ will serve for a $\delta$-net for $B_1$. The relative compactness of $\widetilde{\pi}_m B_1$ follows from the fact that, by Lemma 2, $\widetilde{\pi}_m B_1$ is a

bounded subset of the finite-dimensional space $\mathscr{L}_m$. Thus $B_1$ is shown to be relatively compact and hence a converging subsequence $\{a_1^{n'}\} \subset B_1$, $a_1^{n'} \to a_1 \in \ell_2$ (in norm), can be separated.

Consider the set $B_2 = \{a_2^{n'}\}$ now. By the same method as used above, a subsequence $\{a_2^{n''}\}$, $a_2^{n''} \to a_2 \in \ell_2$, can be isolated. Then it becomes clear that the subsequence $\{A_{n''}\}$ converges, in Hausdorf metrics, to the set $A = \{a_1, a_2\}$,

$$h(A_{n''}, A) \to 0, \ n'' \to \infty. \tag{8.3}$$

Finally, show that $A \in \mathscr{A}_2^*(P)$. Due to the continuity property of the mapping $A \to W(A, P)$ (see Appendix, Lemma A2) from (8.3) it follows that

$$W(A_{n''}, P) \to W(A, P), \ n'' \to \infty.$$

At the same time, $\{A_{n''}\}$ is a minimizing sequence and thus

$$W(A_{n''}, P) \to W_k(P).$$

Hence $W(A, P) = W_k(P)$ and so $A \in \mathscr{A}_2^*(P)$. This proves the theorem.

Theorem 3 enables us to obtain a further result.

Theorem 4. Assume that $P$ is a measure of the second order in $\ell_2$ and the sequence $\{P_n\}$ is such that

a) $P_n \Rightarrow P$,

b) $\|x\|^2$ is uniformly integrable with respect to $P_n$, $n \geqslant 1$.

Then, for any sequence $\{A_n^*\}$ satisfying $A_n^* \in \mathscr{A}_2^*(P)$,

$$h(A_n^*, \mathscr{A}_2^*(P)) \to 0, \ n \to \infty. \tag{8.4}$$

Proof. Due to Corollary 1 we have $W(A_n^*, P) \to W_2(P)$. Hence $\{A_n^*\}$ is a minimizing sequence and the assertion 2) of Theorem 3 may be applied to get (8.4).

Remark. In the case when $\mathscr{A}_2^*(P)$ consists of a single set $A^*$ the relation (8.4) reduces to the simple convergence

$$h(A_n^*, A^*) \to 0.$$

If $P_n$ is the empirical measure, then, by Lemma 5, the latter theorem is valid with probability 1. Such a result may be interpreted as a generalization of the strong law of large numbers (SLLN). In our terms, the SLLN asserts the almost sure convergence of the sequence of the empirical 1-centres $a_n^* = (x_1 + \ldots + x_n)/n$ to the population 1-centre $a^* = EX$.

## 9. Appendix.

Two auxiliary results will be given here. At first, let us define

$$d(x,A) = \inf_{a \in A} d(x,a).$$

**Lemma A1.** If $\int \varphi(d(x,y_0))P(dx) < \infty$ for some $y_0 \in T$, then, for any $A \subset T$ $(A \neq \emptyset)$ and any $R \geqslant 0$, also

$$\int_T \varphi\,(d(x,A) + R)P(dx) < \infty \ .$$

**Proof.** Since $d(x,A) \leqslant d(x,a)$, for all $a \in A$, the monotony property of $\varphi$ implies that

$$\int_T \varphi(d(x,A)+R)P(dx) \leqslant \int_T \varphi(d(x,a) + R)P(dx).$$

Show the latter integral is finite. Let $B = B(a,R)$, $\bar{B} = T \setminus B$. Then, for the monotony and $\Delta_2$-property of $\varphi$, we have

$$\int_T \varphi(d(x,a)+R)P(dx) \leqslant \int_B \varphi(2R)P(dx) + \int_{\bar{B}} \varphi(2d(x,a))P(dx) \leqslant$$

$$\leqslant \varphi(2R) + \lambda \int_{\bar{B}} \varphi(d(x,a))P(dx).$$

Further, let $B_1 = \{x : d(x,y_0) < d(a,y_0)\}$. Then, by the triangle inequality, it is seen that

$$\int_{\bar{B}} \varphi\,(d(x,a))P(dx) \leqslant$$

$$\leqslant \int_{B_1} \varphi(d(x,y_0)+d(a,y_0))P(dx) + \int_{\bar{B}_1} \varphi(d(x,y_0)+d(a,y_0))P(dx).$$

In view of $\Delta_2$-property, the integral over $B_1$ does not exceed the quantity $\lambda \varphi(d(a,y_0))\, P(B_1) < \infty$ and the other integral does not exceed the integral $\lambda \int_T \varphi(d(x,y_0))P(dx)$, which is finite, by our assumption. This proves the lemma.

In the following lemma the sets A and $A_n (n \geqslant 1)$ will be allowed to be arbitrary subsets of a separable metric space T. Let us define

$$W(A,P) = \int_T \varphi\,(d(x,A))P(dx)\ , \quad A \subset T,\ A \neq \emptyset.$$

It is clearly seen that in the case when $|A| = k$ this definition reduces to (1.1).

**Lemma A2.** If $\int \varphi(d(x,y_0))P(dx) < \infty$ for some $y_0 \in T$, then from $h(A_n,A) \to 0$ it follows that

$$W(A_n,P) \to W(A,P). \tag{A.1}$$

**Proof.** It is an easy exercise to show that for any A and B the inequality

$$|d(x,A) - d(x,B)| \leqslant h(A,B) \tag{A.2}$$

holds. Hence, by continuity of $\varphi$, the convergence $h(A,A_n) \to 0$ implies $\varphi(d(x,A)) \to \varphi(d(x,A_n))$ for any $x \in T$. To get (A.1) we apply Lebesgue's theorem. According to this theorem, it suffices to point out an integrable function $g(x)$ which majorizes all the functions $\varphi(d(x,A_n))$ for n greater than

some N. Show g(x) may be taken as $g(x) = \psi[d(x,A)+1]$, for example. Indeed, by $h(A,A_n) \to 0$, $h(A_n,A) < 1$ for every n larger than some N. In view of (A.2), this implies that

$$d(x,A_n) \leq d(x,A)+1, \quad n \geq N.$$

Furthermore, by monotony of $\psi$ ,

$$\psi(d(x,A_n)) \leq \psi[d(x,A)+1] \equiv g(x) , \quad n \geq N,$$

i.e., g(x) majorizes all $\psi(d(x,A_n))$ with n large sufficient. The integrability of g(x) follows directly from Lemma A1.

This completes the proof.

## References

1. Abaya, E.F., Wise, G.L. Convergence of vector quantizers with applications to optimal quantization. SIAM J.Appl. Math., 1984, 44, 183-189.

2. Pärna, K. Strong consistency of k-means clustering criterion in separable metric spaces. Acta et Commentationes Universitatis Tartuensis, 1986, 733, 86-96.

3. Pollard, D. Strong consistency of k-means clustering. Ann. Statist., 1981, 9, 135-140.

4. Ranga Rao, R. Relations between weak and uniform convergence of measures with applications. Ann. Math. Statist., 1962, 33, 659-680.

5. Varadarajan, V.S. On the convergence of sample probability distributions. Sankhya, 1958, 19, 23-26.

6. Биллингсли П. Сходимость вероятностных мер. М., 1977

7. Люстерник Л.А., Соболев В.И. Краткий курс функционального анализа. М., 1982

ОБ УСТОЙЧИВОСТИ МЕТОДА к-СРЕДНИХ В МЕТРИЧЕСКИХ
ПРОСТРАНСТВАХ

К.Пярна

Р е з ю м е

Пусть P, $P_n$ ($n \geq 1$)-вероятностные меры на сепарабельном метрическом пространстве (T, d). Доказано, что $\inf_A W(A,P_n) \to \inf_A W(A,P)$, где функционал $W(\cdot,\cdot)$ определен в (I.I). В пространстве $l_2$ доказана и сходимость $P_n$-оптимальных множеств $A_n^*$. Статья продолжает работу [2].

# A REPRESENTATION OF CLUSTER-ANALYSIS: THE APPROACH
## USING MONOTONIC SYSTEMS
### R.Ääremaa, K.Ääremaa, T.Tamman

The possibility to construct the specific monotonic systems' kernels which will be identical to clusters formed by some well-known cluster-methods was demonstrated in [1]. The methods chosen for demonstration in [1] have the common quality - the result of clustering is based on the ordering of the values of similarities between the elements to be clustered and is not based on the computed real values of similarities. In this paper we set the goal to demonstrate the realization of cluster-analysis, where 1) the values of similarities will be considered in the process of clustering (not only their ordering will be used); 2) no requirements about the shape, or number, or degree of intersection of the clusters will be postulated. One of the possibilities to carry out such a clustering is to construct the monotonic systems without restrictions used in [1]. In the discussion of the problems of clustering we use graph-representation.

## 1. Graph-representation and comparison of the subgraphs

Let $X = \{x_1, x_2, \ldots, x_n\}$ be the set of the elements to be clustered and $S = (s_{ij})$, $i,j=1,\ldots,n$, - the matrix of similarities computed on $X$, where $s_{ij}$ corresponds to the similarity between the elements $x_i$ and $x_j$; $s_{ij} = s(x_i, x_j)$.

The pair $G = (X, E^S)$ presents the graph, where the vertices are defined by $X$ and the real value $s_{kl}$ defined by $S$, corresponds to the edge $(x_k, x_l)$ of the set of the edges $E^S = \{(x_i, x_j): i \neq j, i,j=1,\ldots,n\}$.

For estimating the mutual connectedness of the vertices in the graph, we define a weight function $\pi$, which assigns to each vertex $x_i$ a non-negative real value $\pi(x_i)$ - the weight of the vertex $x_i$. For example the weight of the vertex $x_i$ may be defined as the mean of the mutual similarities between $x_i$ and other vertices by the formula

37

$$\pi(x_i) = \frac{1}{n-1} \sum_{j=1, j\neq i}^{n} s_{ij} . \qquad (1)$$

Let us denote a subgraph of graph G by G', where $G' = (X', E')$ and $X' \subseteq X$, $E' \subseteq E^S$.

We want to search for the subgraph $\bar{G} = (\bar{X}, \bar{E})$, which is possible to consider to be "the best" among all subgraphs of G, taking into consideration the weights of the vertices of each subgraph. We cannot inspect the subgraphs and make comparisons between them to find "the best", presupposing them to be separate graphs. The found subgraph is to rise to the fore against the background of the graph G or of some other graph which has taken to be the base-graph.

Let us take the base-graph for the graph G the graph $G_o = (X, E^o)$, which can be interpreted as a graph, where the mutual similarities between n elements to be clustered are not made known and we consider each of them to be equal to zero. The process of transition from $G_o$ to G is feasible by characterizing the edges of $G_o$ using the values of matrix S. To make comparisons between different subgraphs $G' = (X', E')$ of G, we observe the changes of the weights of the vertices which arise in $G_o$ after estimating the edges by the values of S.

Let us use the following criterion: the subgraph $\bar{G}$ of G is "the best", if the minimum change of the weights of its vertices has the maximum value, where maximization is taken over all subgraphs of G. This subgraph $\bar{G}$ will be called the kernel of the graph G.

**Example 1.** Let $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ and the graph $G = (X, E^S)$ is presented in Figure 1, where the edges are characterized by the values of similarities and the vertices by their weights computed by the formula (1).

Let us choose two subgraphs

$G^1 = (\{x_1, x_2, x_3\}, \{(x_1 x_2), (x_1 x_3), (x_2 x_3)\})$ and

$G^2 = (\{x_1, x_2, x_3, x_4\}, \{(x_1 x_2), (x_1 x_3), (x_1 x_4), (x_2 x_3), (x_2 x_4), (x_3 x_4)\})$

for the comparison. Picturing these subgraphs on base-graph $G_o$ (Figure 2) we can estimate the changes induced on it. We can see that

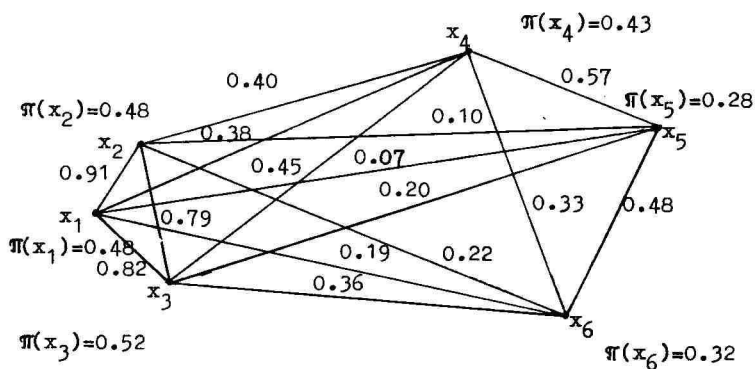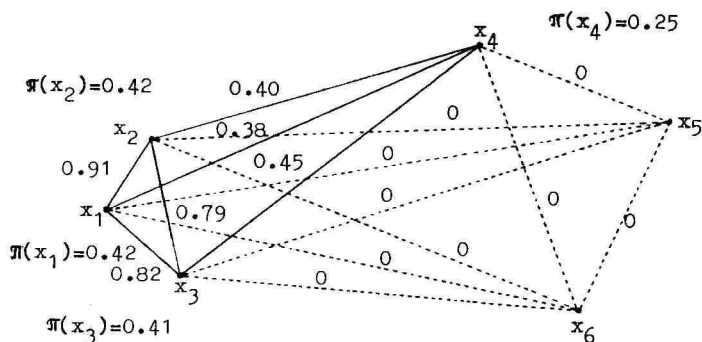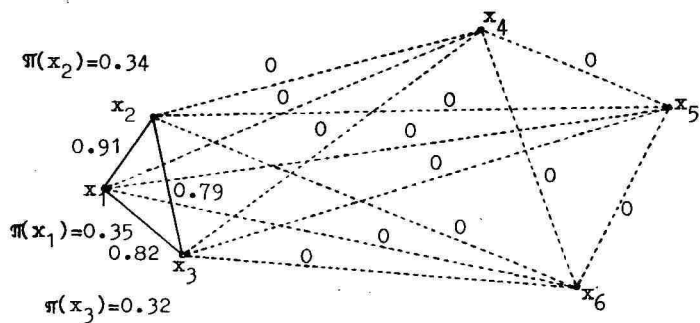$$\min_{x \in G^1} \pi(x) = 0.32 > 0.25 = \min_{x \in G^2} \pi(x).$$

Figure 1



Figure 2

On the basis of the criterion used by us, the subgraph $G^1$ is "better" than $G^2$. It is not difficult to control that there does not exist any subgraph of G by which the minimal change of the weights of its vertices is greater than 0.32. So $G^1$ is the kernel of the graph G in the above given sense.

In general the kernel of the graph depends on the used weight function. The weight function must express the degree of mutual connectedness of elements to be clustered - thus the weight function depends on the measure of the connection (similarity function) to be chosen.

The concept of the kernel of the graph can be modified and defined in such a way that only the edges, which have the value equal or greater than the received minimum value of changes of the weights of vertices, are considered. In such a case the kernel must not be a complete subgraph - it can even consist of several disconnected parts.

Till now we have dealt with the problem of finding only one, the most important subgraph - the first kernel. The whole process of searching of the kernel can be repeated on the subgraph, the elements of which do not belong to the first kernel. There are several different possibilities to eliminate the influence of the elements of the kernel to the whole graph, three of them will be produced below to construct the monotonic systems for clustering.

The kernel $\overline{G} = (\overline{X}, \overline{E})$ was found with respect to the graph $G_o = (X, E^o)$. The graph $\overline{G}$ can be considered to be an independent graph and it is possible to search the kernel of this graph with the respect to the base-graph $\overline{G}_o = (\overline{X}, \overline{E}^o)$. The transition from $\overline{G}_o$ to $\overline{G}$ may be realized using the values of the mutual similarities between the vertices from $\overline{X}$. So a recursive process for finding the kernels in the kernels can be carried out.

All the process described above can be modelled using the theory of monotonic systems.

## 2. Constructing a monotonic system for cluster-analysis

The pair $\langle W, \pi \rangle$ is called to be the monotonic system on the set of elements W, if for any element $a \in W'$, $W' \subseteq W$, the weight function $\pi(a, W')$ satisfies the condition of monotony

$$\pi(a, W'') \leqslant \pi(a, W')$$

40

for any subset $W'' \subseteq W' \subseteq W$.

The set of elements $\overline{W} \subseteq W$ is the kernel of the monotonic system $\langle W, \pi \rangle$, if

$$F(\overline{W}) = \max_{W' \subseteq W} \ (\min_{a \in W'} \ \pi(a, W')).$$

There exists an algorithm (see [1]), which in finite number of steps produces the kernel of the monotonic system. It is possible to define the concrete monotonic system in such a way that its kernel coincides with the kernel of the graph described above. But it is not our aim. We used the graph-representation only for demonstration how the kernel of the monotonic system constructed for clustering can be interpreted using the graphs. If the weight function of the monotonic system expresses the mutual connectedness (simi-larity) of the elements, the kernel may be treated to be a cluster (see [3]).

Let us now construct a concrete monotonic system for clustering the objects $X = \{x_1, x_2, \ldots, x_n\}$. To realize the possibility to get intersecting clusters (overlapping clus-tering) the monotonic system $\langle W, \pi \rangle$ is defined so that the set $W$ is a union of two sets $X$ and $E$, where $X$ is the set of the objects and $E$ is the set of the pairs of these objects (the ties between the objects)
$E = \{(x_i, x_j): \ i \neq j, \ i, j = 1, \ldots, n\}$, where $x_i, x_j \in X$ and $(x_i, x_j) = (x_j, x_i)$. Hence, $W = X \cup E$.

The weight of the object $x_i \in X$ may be defined in several ways. It may be equal to the mean of the mutual similarities between $x_i$ and other objects $x_j \in X$, $x_i \neq x_j$, as in the example given above. It may also be equal to the root-mean-square of the mutual similarities between $x_i$ and other objects $x_j \in X$, $x_i \neq x_j$ as in [3] and in the example given below. In general it is possible to define the weight function in a different manner, but in any case it is necessary to observe that the monotony of the system should be guaranteed.

It is easy to prove that the monotony of the system $\langle X \cup E, \pi \rangle$ is guaranteed if the weight function $\pi$ is defined on any subset $W' = X' \cup E'$ for ties as:

$$\pi((x_i, x_j), E') = \begin{cases} s(x_i, x_j), & \text{if } x_i, x_j \in X' \\ 0 & \text{in other cases} \end{cases} \tag{2}$$

41

and for objects, using natural number  p  (p≠0),  as:

$$\pi(x_i,X') = \begin{cases} (\dfrac{1}{n-1}\displaystyle\sum_{(x_i,x_j)\in E'} s^p(x_i,x_j))^{\frac{1}{p}}, \text{if } \exists x_j \in X':(x_i,x_j)\in E' \\ \\ 0 \qquad\qquad\qquad , \text{if } \forall x_j \in X':(x_i,x_j)\notin E'. \end{cases} \qquad (3)$$

Let us note  that according to  the given function  the
weight of the tie  $\pi((x_i,x_j),E')$  can have the value greater
than  zero  only  if  both  $x_i, x_j \in X'$.  Let us call  the tie
$(x_i,x_j)$  to be  single in  X'  if at least  $x_i$  or  $x_j$  does
not belong to  X'. The single object is  such an object which
has no ties.  It is evident  that the kernel  of this system
⟨X∪E,π⟩ has neither any single tie  nor any single object.

Example 2. Let us use the given monotonic system for clus-
tering the data presented  by the matrix of similarities  S,
values of which were indicated in Figure 1.  We can find the
first cluster as the kernel of the defined monotonic system,
where  $W = \{x_1,x_2,x_3,x_4,x_5,x_6,\ (x_1x_2),(x_1x_3),\dots,(x_5x_6)\}$  and
the weight function  $\pi$  is defined for ties by (2)  and  for
objects by (3) using value  p=2. The set  W  and the weights
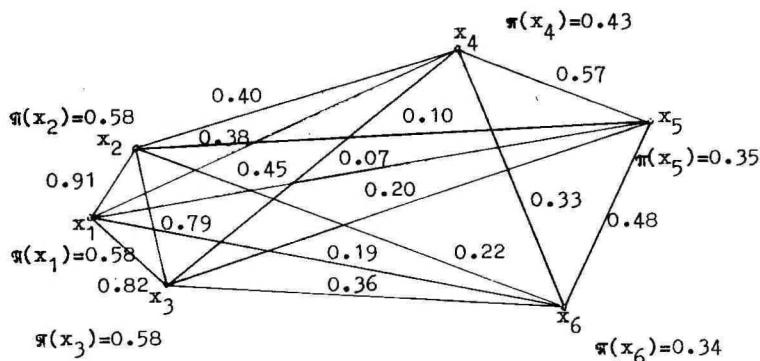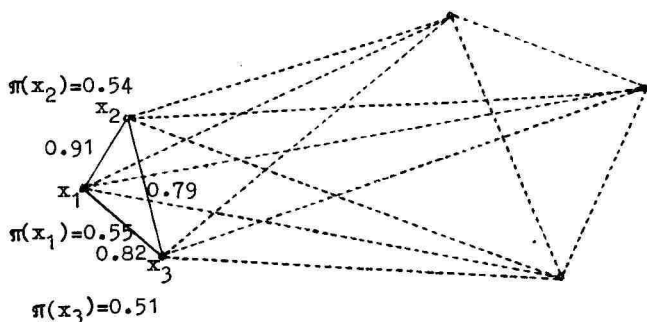of its elements are presented in Figure 3.



Figure 3

42

$\pi(x_2)=0.54$

$x_2$

$0.91$

$x_1$ $0.79$

$\pi(x_1)=0.55$

$0.82$

$x_3$

$\pi(x_3)=0.51$

Figure 4

Using the algorithm for finding the kernel, we can estimate that the set
$$\overline{W} = \overline{X} \cup \overline{E} = \{x_1, x_2, x_3, (x_1x_2), (x_1x_3), (x_2x_3)\}$$
is the kernel, its structure is presented in Figure 4.

### 3. Constructing the monotonic systems for non-overlapping or overlapping clusters

It is possible to construct the sequence of the monotonic systems for searching either non-intersecting or intersecting clusters. Let us have the sequence of l systems
$$\langle W^0, \pi \rangle, \quad \langle W^1, \pi \rangle, \ldots, \langle W^1, \pi \rangle,$$
where $W^i = X^i \cup E^i$, $i=0,1,\ldots,l$, and $W^0 \equiv W$, $X^0 \equiv X$, $E^0 \equiv E$, and $\overline{W}^i = \overline{X}^i \cup \overline{E}^i$ denotes the kernel of the i-th monotonic system.

To get the non-overlapping clustering, which is made up of l non-intersecting clusters, the k-th system of the sequence of l systems is defined eliminating all elements (objects and ties) of the found kernel $\overline{W}^{k-1}$ from the (k-1)-th system. The elements which will turn into single after such a removal are also omitted. So it is possible to define the k-th monotonic system, $k=1,2,\ldots,l$, as

$\langle W^k, \pi \rangle$, where $W^k = X^k \cup E^k$ and
$$X^k = X^{k-1} \setminus (\overline{X}^{k-1} \cup \hat{X}^{k-1}), \quad E^k = E^{k-1} \setminus \hat{E}^{k-1},$$
$$\hat{X}^{k-1} = \{a: (a,b) \notin E^k \text{ for any } b \in X^{k-1}\},$$
$$\hat{E}^{k-1} = \{(a,b): a \in \overline{X}^{k-1}\}.$$

43

The number 1 is defined by $X^{1-1} = \emptyset$.

To get the overlapping clustering of 1 clusters, we define the sequence of 1 monotonic systems in which each next system is received by omitting either the ties belonging to the kernel of the previous system or the ties between the objects of the found kernel but not omitting the objects of the kernel. In any case the elements which become single will be omitted too.

The first of the above mentioned two possibilities to receive the k-th monotonic system on the basis of (k-1)-th system, $k = 1,2,\ldots,1$, is produced so: for the k-th monotonic system $\langle W^k, \pi \rangle$, where $W^k = X^k \cup E^k$,

$$X^k = X^{k-1} \setminus \hat{X}^{k-1}, \quad E^k = E^{k-1} \setminus \bar{E}^{k-1},$$

$$\hat{X}^{k-1} = \{a: \ (a,b) \notin E^k \ \text{for any} \ b \in X^{k-1}\}.$$

The second possibility to receive the k-th monotonic system can be presented in such a way: for the k-th monotonic system $\langle W^k, \pi \rangle$, where $W^k = X^k \cup E^k$,

$$X^k = X^{k-1} \setminus \check{X}^{k-1}, \quad E^k = E^{k-1} \setminus \check{E}^{k-1},$$

$$\check{X}^{k-1} = \{a: \ (a,b) \notin E^k \ \text{for any} \ b \in X^{k-1}\},$$

$$\check{E}^{k-1} = \{(a,b): \ a,b \in \bar{X}^{k-1}\}.$$

**Example 3.** In the previous example we found the set

$$\bar{W} \equiv \bar{W}^o \equiv \bar{X}^o \cup \bar{E}^o = \{x_1, x_2, x_3, (x_1x_2), (x_1x_3), (x_2x_3)\}$$

to be the kernel of the first monotonic system in the sequence of the systems.

To construct the overlapping clusters, we define the next system $\langle W^1, \pi \rangle$ omitting the ties of the found kernel $\bar{W}^o$ from the preliminary system $\langle W^o, \pi \rangle$. There is not any element which would become single after such a removal. Hence, we received the set $W^1$, which includes the following objects: $x_1$, $x_2$, $x_3$, $x_4$, $x_5$, $x_6$
and the ties: $(x_1x_4)$, $(x_1x_5)$, $(x_1x_6)$, $(x_2x_4)$, $(x_2x_5)$, $(x_2x_6)$, $(x_3x_4)$, $(x_3x_5)$, $(x_3x_6)$, $(x_4x_5)$, $(x_4x_6)$, $(x_5x_6)$.

The set $W^1$ with the weights of its elements is presented in Figure 5.

It is possible to find the kernel of the system $\langle W^1, \pi \rangle$ to be the set $\bar{W}^1 = \{x_4, x_5, x_6, (x_4x_5), (x_4x_6), (x_5x_6)\}$ and the kernel of the next system $\langle W^2, \pi \rangle$ as the set $\bar{W}^2 = \{x_3, x_4, (x_3x_4)\}$, and so on.

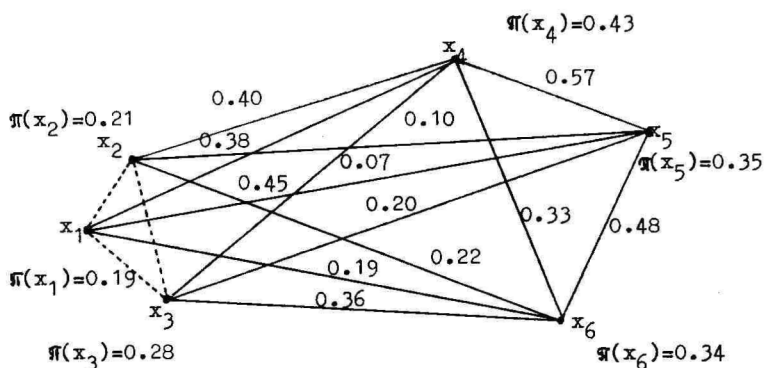Figure 5

## 4. Clustering of objects and attributes, and simultaneous clustering of objects and attributes

The similar methodology which is based on the monotonic systems as described for clustering of the objects X can be used for clustering of the attributes Y. It is possible to construct a monotonic system on the attributes and the ties between them. The weight function for the attributes may be defined on the basis of the used similarity function computed on the attributes and it expresses the mutual connection of the attributes to be clustered.

The objects and the attributes can be treated to be the elements of the same monotonic system or to be the elements of two related systems, and so the mutual effect of the objects and the attributes can be examined. In a simpler case the problem can be raised how to search for these quite strongly connected attributes which correspond to some cluster consisting of quite strongly connected objects. To solve this problem the monotonic system on the attributes and their ties is constructed, where the attributes are presented only by objects belonging to the observed cluster. Such an approach is useful in some kinds of data analysis, but in fact it is not simultaneous clustering of objects and attributes.

The simultaneous clustering must reveal the intrinsic interaction structure of objects and attributes (show which groups of the objects and attributes interact together). For the simultaneous clustering of the set of objects and attributes $Z = X \cup Y$ the monotonic system is constructed

45

regarding the objects with their ties and the attributes with their ties as elements, equal in rights, of the monotonic system and the weight function is defined so that the values of the weights of all the elements (objects, ties between objects, attributes, ties between attributes) are located in the same interval (for example in $[0,1]$).

It is necessary to define the weight function so that the monotony of the system should be guaranteed. We have defined the weight function on the matrix of similarities S between the objects and on the matrix of similarities T between the attributes. In order to guarantee the monotony of the constructed system, the following condition must be satisfied: the removal of any attribute from the set Z may give rise to the change of the values of the matrix S only in the direction of decreasing and the removal of any object from Z may give rise to the change of the values of the matrix T only in the direction of decreasing. Hence, the similarity functions defined on the objects and on the attributes, which give the values as the matrix S and T accordingly, must be in harmony required by this condition.

<u>Example 4.</u> Let us present the data by the matrix

$$\begin{pmatrix} 1.0 & 1.0 & 3.0 & 5.0 & 7.0 & 9.0 \\ 0.5 & 1.5 & 1.0 & 3.0 & 6.0 & 5.5 \\ 0.9 & 0.7 & 1.8 & 2.0 & 4.2 & 4.0 \\ 2.0 & 2.0 & 2.0 & 1.0 & 1.0 & 2.0 \end{pmatrix}$$

where the rows characterize the attributes $Y = \{y_1, y_2, y_3, y_4\}$ and the columns - the objects $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$.

Let the similarity $s_{ij}$ between the objects $x_i$ and $x_j$ be computed using the formula

$$s_{ij} = 1 - (\frac{1}{m} \sum_{k=1}^{m} (x_i^{(k)} - x_j^{(k)})^2)^{\frac{1}{2}},$$

at which $x_i^{(k)}$ represents the value of the k-th attribute of the i-th object after standardizing every single one out of m=4 attributes in the interval $[0,1]$.

We shall regard the simplest case, where the similarity $t_{ij}$ between the attributes $y_i$ and $y_j$ is computed using the similar formula as above

$$t_{ij} = 1 - (\frac{1}{n} \sum_{k=1}^{n} (y_i^{(k)} - y_j^{(k)})^2)^{\frac{1}{2}},$$

at which $y_i^{(k)}$ represents the value of the k-th object for the i-th attribute standardized in the interval $[0,1]$. The values of the computed matrixes S and T are indicated in Figure 6 as the weights of the ties between the objects and between the attributes accordingly. Let us note that the computed similarities coincide with the similarities presented in Figure 3.
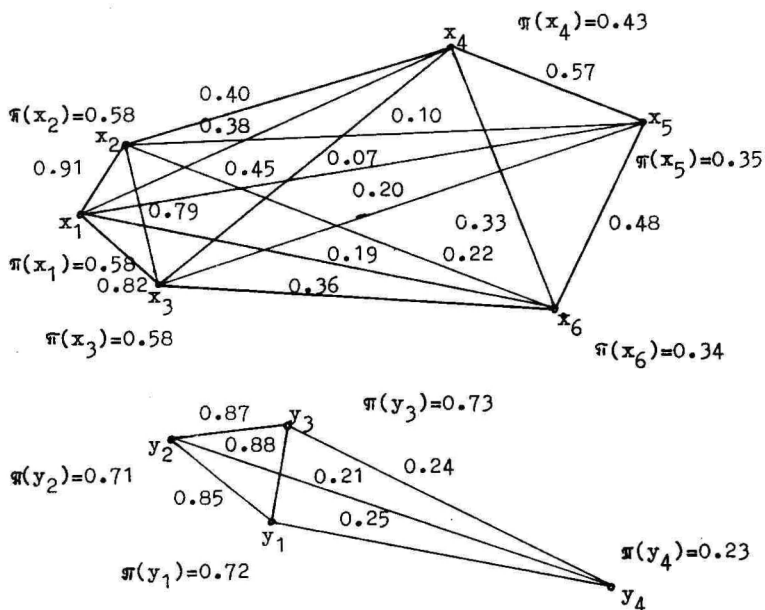


Figure 6

The found kernel includes the elements
$x_1, x_2, x_3, x_4$, $(x_1 x_2)$, $(x_1 x_3)$, $(x_1 x_4)$, $(x_2 x_3)$, $(x_2 x_4)$, $(x_3 x_4)$, $y_1, y_2, y_3$, $(y_1 y_2)$, $(y_1 y_3)$, $(y_2 y_3)$ - it is presented in Figure 7 by help of two graphs.

Let us use the construction of the monotonic systems for clustering objects and attributes allowing to overlap the objects between themselves and the attributes between themselves in the clusters.

The constructed system, which includes the elements
$x_1, x_2, x_3, x_4, x_5, x_6$, $(x_1 x_5)$, $(x_1 x_6)$, $(x_2 x_5)$, $(x_2 x_6)$, $(x_3 x_5)$, $(x_3 x_6)$, $(x_4 x_5)$, $(x_4 x_6)$, $(x_5 x_6)$,
$y_1, y_2, y_3, y_4$, $(y_1 y_4)$, $(y_2 y_4)$, $(y_3 y_4)$,
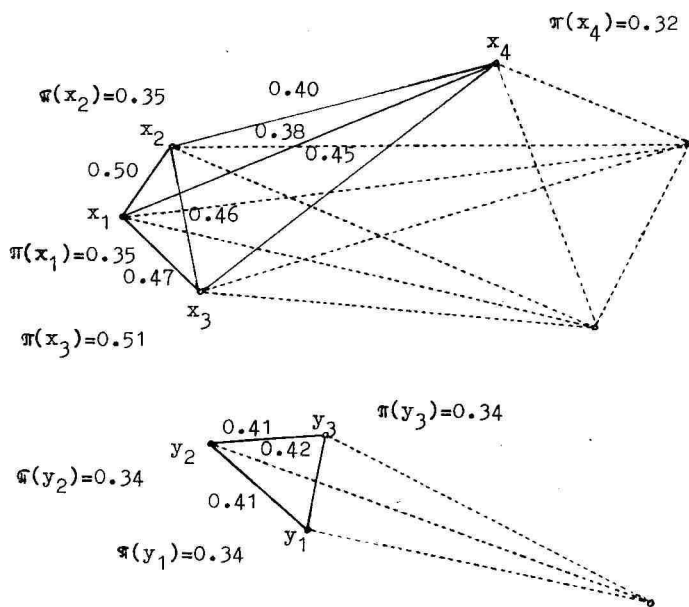is presented in Figure 8 using two graphs.
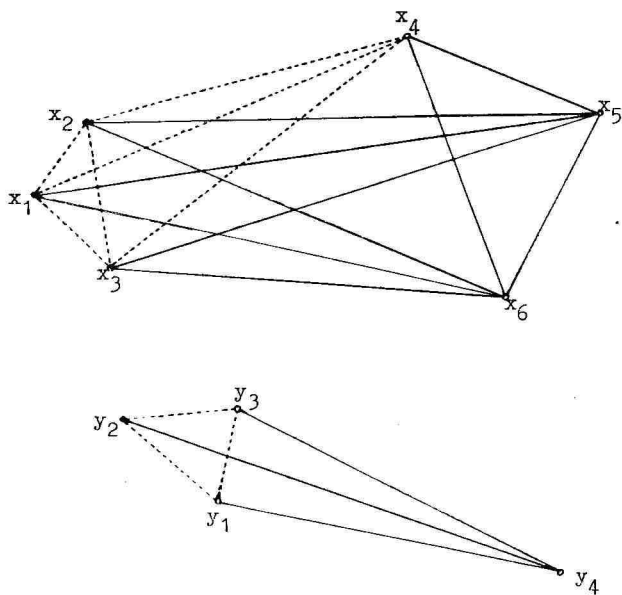
47

Figure 7



Figure 8

49

The kernel of the corresponding system includes the elements
$x_3, x_4, x_5, x_6,$ $(x_3 x_5),$ $(x_3 x_6),$ $(x_4 x_5),$ $(x_4 x_6),$ $(x_5 x_6)$ and
$y_1, y_3, y_4,$ $(y_1 y_4),$ $(y_3 y_4).$

## 5. Some aspects of the cluster-analysis using the monotonic systems

In the previous treatment we did not bring forth the values of the weights of the elements of the found kernels. The minimum of the weights of the elements of the kernel of the monotonic system, constructed by omitting the kernel already found completely or partially from the previous system, is less than the corresponding minimum computed for the kernel of the previous system. The minimum of the weights of its elements characterizes each found kernel. This value may be observed as the value of the level of the level-clustering and the whole process of searching the clusters can be regarded as constructing stratified clustering (see [2,3]). At that clustering the non-intersecting clusters of the level-clustering as the subsets of the kernel, which have no ties between their elements may be found; of course each subset itself includes the ties.

The full treatment presented above may be used for any found kernel (cluster) and an independent monotonic system may be constructed on the elements of this kernel (cluster).

Thus, it is possible to find a specific structure of clusters with different degrees of connectedness - these clusters may be overlapped, and at the same time it is possible to point out the set of the most strongly connected elements as a cluster in each cluster in turn. Which elements of the cluster are strongly connected, which weakly and which are not connected at all may be explained, because any tie belonging to the cluster is characterized by a real value of its weight.

The cluster-analysis carried out using the monotonic systems makes possible to simplify the data by giving a structure of the elements which are under investigation, presenting no requirements about the shape, or number, or degree of intersection of the clusters, and to examine the received structure more precisely if such a exploration is required.

References

1. Aaremaa R., Clustering as a process of finding of kernels of monotonic systems. Acta et Commentationes Universitatis Tartuensis, 1984, 685, 64-75.

2. Aaremaa R., Aaremaa K., On the method of searching stratified clustering. Short communications and posters, Rome, 1986.

3. Ээремаа Р., Конструирование послойной кластеризации. Труды ВЦ ТГУ, 1986, 53, II7-I33.

ПРЕДСТАВЛЕНИЕ КЛАСТЕРНОГО АНАЛИЗА:
ПОДХОД НА ОСНОВЕ МОНОТОННЫХ СИСТЕМ
Р.Ээремаа, К.Ээремаа, Т.Тамман
Р е з ю м е

Проблема кластеризирования интерпретируется проблемой выделения наиболее существенных подграфов графа близости. При этом значимость каждого подграфа оценивается на основе его сопоставления с некоторым выбираемым базой графом. Описываемый подход хорошо согласуется с развиваемой авторами методикой кластеризирования данных с помощью монотонных систем, является интерпретацией этой методики. В зависимости от определения конкретной монотонной системы, т.е. от выбора элементов системы и весовой функции, получаются системы кластеризирования либо объектов, либо признаков, либо объектов и признаков одновременно. Приводится необходимое условие для конструирования систем проведения одновременного кластеризирования объектов и признаков. Включение в элементы системы связей между кластеризируемыми единицами дает возможность определить кластеры с учётом существенных связей, а также конструировать кластеризацию с пересекающимися кластерами.

Весь процесс кластеризирования рассматривается как нахождение ядер конструированной специальным образом последовательности монотонных систем. В зависимости от определения этой последовательности получается либо непересекающаяся либо пересекающаяся кластеризация. При этом в пересекающейся кластеризации степень пересекаемости определяется данными; накладывать какие-либо предварительные требования на степень пересекаемости не допускается. Аналогичный кластерный анализ можно провести для элементов уже найденных кластеров.

# DISTRIBUTION APPROXIMATION BY NORMAL MIXTURES
## AND EDGEWORTH EXPANSIONS

T. Kollo    I. Traat    J. Vilismäe

The central limit theorem is the base in constructing the approximate interval estimations for statistics in very many cases. The estimators are obtained by using the limit distributions or different power series for the distribution of regarded statistic. During the recent years many investigators have paid attention to the Edgeworth expansions. Let us propose that for a p-dimensional sample $X = (X_1, \ldots, X_n)$ from a general population with first finite moments statistic $T_n$ admits the Edgeworth expansion

$$F(x) = P(T_n < x) = \Psi(x) + \frac{1}{\sqrt{n}} h(x) \, \varphi(x) + \mathcal{O}(\frac{1}{\sqrt{n}}), \quad (1)$$

where $\Psi(x)$ and $\varphi(x)$ are the distribution and density functions of $N(0,d)$, d is the limiting variance of $T_n$, and $h(x)$ — the function, depending on the first three cumulants of $T_n$ (the expansion holds if $T_n$ is the function of sample mean or of covariance matrix for example).

If statistic $T_n$ is a skewed random variable, it is natural to use a non symmetric distribution for its approximation. Hall [1] considered $T_n$ as the sum of independent identically distributed random variables with zero mean and unit variance and got the approximation for $F(x)$ by $\chi^2$-distribution. For fitting the best $\chi^2$-distribution, the third moments of summands were used. The accuracy of the approximation was the same or even better than using two first terms in the expansion (1). We are interested in approximating the distribution function $F(x)$ with the mixture of normal distributions that permits to treat multivariate and univariate statistics by the same method.

Let a random variable Z be the mixture of $X \sim N(\mu_1, \delta^2)$ and $Y \sim N(\mu_2, \delta^2)$ with the distribution function

$$F_Z(x) = \gamma F_X(x) + (1 - \gamma) F_Y(x),$$

where $0 \leqslant \gamma \leqslant 1$. Through the standard normal distribution we have

$$F_Z(x) = \gamma \Phi(\frac{x - \mu_1}{\delta}) + (1 - \gamma) \Phi(\frac{x - \mu_1}{\delta} + \frac{\mu_1 - \mu_2}{\delta}). \quad (2)$$

The parameters $\gamma$, $\mu_1$, $\mu_2$ and $\delta^2$ we find with the help of known cumulants $\varpi_k (k = 1, \ldots, 4)$ of random variable Z from the system of equations

$$\varpi_1 = \gamma (\mu_1 - \mu_2) + \mu_2 ; \tag{3}$$

$$\varpi_2 = \gamma (1 - \gamma)(\mu_1 - \mu_2)^2 + \delta^2 ; \tag{4}$$

$$\varpi_3 = \gamma (1 - \gamma)(1 - 2\gamma)(\mu_1 - \mu_2)^3 ; \tag{5}$$

$$\varpi_4 = 6\gamma(1 - \gamma)(\tfrac{3+\sqrt{3}}{6} - \gamma)(\tfrac{3-\sqrt{3}}{6} - \gamma)(\mu_1 - \mu_2)^4 \tag{6}$$

If Z is normal ($\varpi_3 = \varpi_4 = 0$), then $\gamma = 0$ or $\gamma = 1$ or $\mu_1 = \mu_2$. If Z is symmetric nonnormal ($\varpi_3 = 0$, $\varpi_4 \neq 0$), then $\gamma = 1/2$. If Z is nonsymmetric with $\varpi_4 = 0$, then $\gamma = (3+\sqrt{3})/6$ or $\gamma = (3 - \sqrt{3})/6$. In general for $\gamma$ we get from (5) and (6) the equation

$$(216+16a)\gamma^6 - (648+48a)\gamma^5 + (756+56a)\gamma^4 - (432+32a)\gamma^3 + (126+9a)\gamma^2 - (18+a)\gamma + 1 = 0, \tag{7}$$

which has at least two solutions in the interval $(0,1)$,

$$a = \varpi_4^3 / \varpi_3^4 .$$

The other parameters we find from (3) – (5):

$$\mu_1 - \mu_2 = \sqrt[3]{\varpi_3 / [\gamma (1 - \gamma)(1 - 2\gamma)]} ;$$

$$\mu_1 = \varpi_1 + (1 - \gamma)(\mu_1 - \mu_2) ; \tag{8}$$

$$\mu_2 = \varpi_1 - \gamma(\mu_1 - \mu_2) ; \tag{9}$$

$$\delta^2 = \varpi_2 - \gamma (1-\gamma)(\mu_1 - \mu_2)^2 . \tag{10}$$

The statistic $T_n$ has the cumulants of following order

$$\varpi_1 = O(n^{-1/2}), \quad \varpi_2 = O(1), \quad \varpi_3 = O(n^{-1/2}), \quad \varpi_4 = O(n^{-1})$$

Then, if $n \to \infty$ ,

$$a = O(n^{-1}) \to 0,$$

$$\gamma \to (3 + \sqrt{3})/6, \quad (1 - \gamma) \to (3 - \sqrt{3})/6$$

and

$$\mu_1 - \mu_2 \to 0.$$

Consequently, the mixture (2) for the statistic $T_n$ also converges to the normal distribution.

In the equations (3) – (6) and in the Edgeworth expansion we have theoretic population cumulants. In everyday data analysis we can use only their empirical estimates. It has not been theoretically investigated how this replacement influences the accuracy of the approximation of the exact distribution function. In this paper we shall look a simple example to compare the exactness of different approximations.

The statistic we shall consider is

$$T_n = \sqrt{n} ( S^2 - 1),$$

where $S^2$ is the sample variance statistic:
$$S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})^2 , \quad \overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i .$$

The sample is drawn from onedimensional normal population $N(0,1)$. Theoretically
$$Y_n = (n - 1) S^2 \sim \chi_{n-1}^2 ,$$
which determines the exact distribution $F(x)$ of $T_n$:
$$F(x) = P(T_n < x) = P(Y_n < (n-1)(\frac{x}{\sqrt{n}} + 1)) .$$

We consider the following approximations of $F(x)$:
$F_N(x)$ - asymptotic normal distribution $(n \to \infty)$ ;
$F_E(x)$ - Edgeworth expansion (including the term $n^{-1/2}$);
$F_M(x)$ - mixture of normal distributions ;
$\widetilde{F}_N(x)$, $\widetilde{F}_E(x)$, $\widetilde{F}_M(x)$, $\widetilde{F}_{MN}(x)$ - sample estimators of $F_N(x)$, $F_E(x)$, $F_M(x)$, correspondingly ($F_M(x)$ is estimated in two different ways - by $\widetilde{F}_M(x)$ and $\widetilde{F}_{MN}(x)$).

The distribution functions $F_N(x)$ and $F_E(x)$ have the following form:
$$F_N(x) = \frac{1}{\sqrt{2\pi k_2}} \int_{-\infty}^{x} e^{-\frac{t^2}{2k_2}} dt ,$$
$$F_E(x) = F_N(x) - \frac{\gamma_3}{6\sqrt{n} \, k_2^2} (x^2 - k_2) F_N'(x),$$

where
$$k_2 = \overline{\mu}_4 - \overline{\mu}_2^2 ,$$
$$\gamma_3 = \overline{\mu}_6 - 3\overline{\mu}_2\overline{\mu}_4 - 6\overline{\mu}_3^2 + 2\overline{\mu}_2^3 .$$

$F_M(x)$ is given by the mixture (2), the parameters of which $\gamma$, $\mu_1$, $\mu_2$, $\delta$ are found with the help of the cumulants of statistic $T_n$ from the equations (7) - (10). The expressions for the cumulants $\varkappa_1(T_n)$ are derived using results in $[3]$(pp 106, 441 - 442) :

$$\varkappa_1(T_n) = 0 ; \tag{11}$$
$$\varkappa_2(T_n) = \overline{\mu}_4 - \frac{n-3}{n-1} \overline{\mu}_2^2 ; \tag{12}$$
$$\varkappa_3(T_n) = \frac{1}{\sqrt{n}} ( \overline{\mu}_6 - 3\overline{\mu}_4\overline{\mu}_2 - 6\overline{\mu}_3^2 + 2\overline{\mu}_2^3) + \theta(\frac{1}{\sqrt{n}}) ; \tag{13}$$
$$\varkappa_4(T_n) = \frac{1}{n} (\overline{\mu}_8 - 4\overline{\mu}_6\overline{\mu}_2 - 24\overline{\mu}_5\overline{\mu}_3 + 45\overline{\mu}_4^2 - 276\overline{\mu}_4\overline{\mu}_2^2 +$$
$$96\overline{\mu}_2\overline{\mu}_3^2 + 426\overline{\mu}_2^4) + \theta(\frac{1}{n}). \tag{14}$$

The quantities $\overline{\mu}_i$ denote the central moments of the population distribution, in our example of the distribution $N(0,1)$. For the normal population $N(0, \overline{\mu}_2)$ we can present the equalities (11) - (14) in a much simpler form:
$$\varkappa_1(T_n) = 0 ; \tag{15}$$

53

$$\varkappa_2(T_n) = \frac{2n}{n-1} \, \bar{\mu}_2^2 \; ; \tag{16}$$

$$\varkappa_3(T_n) = \frac{8n\sqrt{n}}{(n-1)^2} \, \bar{\mu}_2^3 \; ; \tag{17}$$

$$\varkappa_4(T_n) = \frac{48 \, n^2}{(n-1)^3} \, \bar{\mu}_2^4 \; . \tag{18}$$

The formulae of sample estimators $\widetilde{F}_N(x)$, $\widetilde{F}_E(x)$, $\widetilde{F}_M(x)$ are received from the formulae of $F_N(x)$, $F_E(x)$, $F_M(x)$, if the theoretic central moments $\bar{\mu}_i$ in them are replaced by the sample moments

$$\bar{m}_i = 1/n \sum_{k=1}^{n} (x_k - \bar{x})^i \; .$$

The function $\widetilde{F}_{MN}(x)$ is the sample estimator of $F_M(x)$, which in this case is determined (instead of the equations (11) - (14)) by the equations (15) - (18), where an unbiased estimate $s^2$ is used for $\bar{\mu}_2$.

The goodness of the approximation is estimated by the distances $d_N$, $d_E$, $d_M$, $\widetilde{d}_N$, $\widetilde{d}_E$, $\widetilde{d}_M$, $\widetilde{d}_{MN}$, each of which is defined as the maximum difference between $F(x)$ and corresponding approximating distribution function over all points in ordinary $\chi^2$-tables [2].

The experiment is carried out as follows. For fixed sample size n the distances $d_N$, $d_M$, $d_E$ are calculated. Then K = 1,000 samples are generated from the distribution N(0,1). For each sample the values of the functions $\widetilde{F}_N(x)$, $\widetilde{F}_E(x)$, $\widetilde{F}_M(x)$, $\widetilde{F}_{MN}(x)$ and the corresponding distances $\widetilde{d}_N$, $\widetilde{d}_E$, $\widetilde{d}_M$, $\widetilde{d}_{MN}$ are calculated. From these values the sample means $E\widetilde{d}_N$, $E\widetilde{d}_E$, $E\widetilde{d}_M$, $E\widetilde{d}_{MN}$ and the standard deviations $\sqrt{D\widetilde{d}_N}$, $\sqrt{D\widetilde{d}_E}$, $\sqrt{D\widetilde{d}_M}$, $\sqrt{D\widetilde{d}_{MN}}$, also the minimum values min $\widetilde{d}_N$, min $\widetilde{d}_E$, min $\widetilde{d}_M$, min $\widetilde{d}_{MN}$ and the maximum values max $\widetilde{d}_N$, max $\widetilde{d}_E$, max $\widetilde{d}_M$, max $\widetilde{d}_{MN}$ are found. All this process is repeated for the sample sizes n = 25, 50, 100, 201, 401. The results are given in Table 1.

From the first three rows in Table 1 we can see that the theoretical Edgeworth expansion is the best approximation to the function $F(x)$, the error of approximation is 0,0047 when n = 25, and it decreases with the rate 1/n.

Mixture $F_M(x)$ is better than limiting normal approximation $F_N(x)$, having smaller approximation errors and converging to $F(x)$ faster than $F_N(x)$. Comparing next four rows 4.-7. with the first ones, we can see that on the average, estimators $\widetilde{F}_N(x)$, $\widetilde{F}_E(x)$, $\widetilde{F}_M(x)$ and $\widetilde{F}_{MN}(x)$ are much worse than corresponding theoretical approximations. The approximation error decreases slowly and is remarkable ($\approx 0.02$)

even while sample size n = 400. It is interesting that all observed sample estimators are almost equal in the sense of the approximation goodness. The average distances and also the standard deviations of the distances (rows 8. - 11.) differ only a little from one to another. As the best estimator we can point out the function $\tilde{F}_{MN}(x)$, but in this case we take into account the additional information about the population distribution. From the last eight rows (12-19) we can see the range in which the distances over all 1000 generated samples vary. It appears that for some samples the approximation error of the estimators $\tilde{F}_N(x)$, $\tilde{F}_E(x)$, $\tilde{F}_M(x)$, $\tilde{F}_{MN}(x)$ can be much larger than in average, especially for small n.

TABLE 1

| n | 25 | 50 | 100 | 201 | 401 |
|---|---|---|---|---|---|
| 1. $d_N$ | 0.0389 | 0.0271 | 0.0190 | 0.0134 | 0.0093 |
| 2. $d_E$ | 0.0047 | 0.0022 | 0.0011 | 0.0005 | 0.0003 |
| 3. $d_M$ | 0.0181 | 0.0131 | 0.0057 | 0.0033 | 0.0019 |
| 4. $E\tilde{d}_N$ | 0.0934 | 0.0618 | 0.0430 | 0.0296 | 0.0209 |
| 5. $E\tilde{d}_E$ | 0.0885 | 0.0572 | 0.0394 | 0.0262 | 0.0184 |
| 6. $E\tilde{d}_M$ | 0.0912 | 0.0596 | 0.0395 | 0.0277 | 0.0188 |
| 7. $E\tilde{d}_{MN}$ | 0.0691 | 0.0456 | 0.0313 | 0.0210 | 0.0142 |
| 8. $\sqrt{D\tilde{d}_N}$ | 0.0523 | 0.0357 | 0.0246 | 0.0172 | 0.0117 |
| 9. $\sqrt{D\tilde{d}_E}$ | 0.0583 | 0.0407 | 0.0288 | 0.0200 | 0.0135 |
| 10. $\sqrt{D\tilde{d}_M}$ | 0.0595 | 0.0415 | 0.0283 | 0.0202 | 0.0136 |
| 11. $\sqrt{D\tilde{d}_{MN}}$ | 0.0484 | 0.0311 | 0.0227 | 0.0149 | 0.0101 |
| 12. min $\tilde{d}_N$ | 0.0376 | 0.0264 | 0.0186 | 0.0133 | 0.0092 |
| 13. min $\tilde{d}_E$ | 0.0046 | 0.0015 | 0.0013 | 0.0010 | 0.0004 |
| 14. min $\tilde{d}_M$ | 0.0090 | 0.0045 | 0.0024 | 0.0014 | 0.0008 |
| 15. min $\tilde{d}_{MN}$ | 0.0135 | 0.0081 | 0.0049 | 0.0028 | 0.0016 |
| 16. max $\tilde{d}_N$ | 0.3261 | 0.2441 | 0.1758 | 0.1304 | 0.0744 |
| 17. max $\tilde{d}_E$ | 0.3278 | 0.2489 | 0.1778 | 0.1330 | 0.0755 |
| 18. max $\tilde{d}_M$ | 0.3244 | 0.2256 | 0.1783 | 0.1136 | 0.0751 |
| 19. max $\tilde{d}_{MN}$ | 0.2679 | 0.1665 | 0.1272 | 0.0919 | 0.0631 |

Comparison of different approximations shows that
in the situation where population distribution and its mo-
ments are unknown the asymptotical normal distribution is
nearly as good as more complicated approximations are. Effect
from using the higher moments in the approximations can be
obtained, if the distribution of population is known. It has
not been investigated how the approximation goodness changes,
if we use unbiased estimates of the moments (which have rigor-
ous form). Some information about the optimal value of k has
been obtained during the experiment. The modelling showed
that it is not necessary to take k as big as 1000.
All the values of calculated variables are nearly the same
for k = 200 already, the tendences and relations between in-
vestigated functions were well remarkable even for k = 100.

References
1. Hall P. Chi squared approximations to the distribution
   of a sum of independent random variables. Ann. Probab.,
   1983, 11, №4, 1028-1036.
2. Большев Л.Н., Смирнов Н.В. Таблицы математической статис-
   тики. М.,1965.
3. Кэндалл М.Дж., Стюарт А. Теория распределений. М., 1965
4. Траат И. Представление неизвестных распределений статис-
   тик с помощью смеси нормальных распределений. Труды ВЦ
   ТГУ, 1984, 51, 126-134.

### АППРОКСИМАЦИЯ РАСПРЕДЕЛЕНИЙ С ПОМОЩЬЮ
### НОРМАЛЬНЫХ СМЕСЕЙ И РАЗЛОЖЕНИЙ ЭДЖВОРТА
Т.Колло, И.Траат, Ю.Вилисмяэ

Р е з ю м е

При помощи статистического моделирования изучается точ-
ность аппроксимации распределения статистики с предельным
нормальным распределением, разложением Эджворта и смесью
нормальных распределений. Точность аппроксимации оценивает-
ся и в случае, где распределение генеральной совокупности
является неизвестным.

# AN ESTIMATE FOR THE VOLUME OF A SET OF POLYGONAL PATTERNS

P.Mikkov

## 1.Introduction

In the present paper we are going to examine the set of patterns, which consist of finite straight lines, located in the unite square, containing T-type junctions only (fig.1, see [1]).
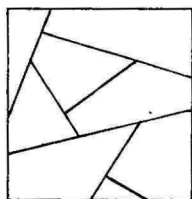


Fig.1.

On the set of these patterns we are going to define the measure $\lambda$ . For this measure we shall describe straight lines, the parts of which make up the pattern, by using of polar coordinates ( see [2]). $\theta$ is the polar angle of the normal line drawn from point O to the straight line, and $\varrho$ is the length of the normal line (fig.2).



Fig.2.

$\Omega$ will denote the set of all patterns, $\Omega_n$ will denote the set of patterns, which consist of n $(n \in \mathbb{N})$ finite straight lines. Therefore

$$\Omega = \bigcup_{n=0}^{\infty} \Omega_n .$$

Let us define measure $\lambda$ on set $\Omega_n$ as

$$d\lambda = d\theta_1 d\varsigma_1 \ldots d\theta_n d\varsigma_n ,$$

conventionally

$$\lambda(\Omega_0) = 1 .$$

Such a measure is related to polygonial Markov fields (see [1]).

The main result of the present paper is the following relation

$$\lambda(\Omega) = \sum_{n=0}^{\infty} \lambda(\Omega_n) < \infty .$$

This relation follows from the theorem, proof of which is given below:

<u>Theorem.</u> The following estimate holds true:

$$\lambda(\Omega_n) < C^n \left( \frac{\ln \ln n}{\ln n} \right)^n$$

where $n \geqslant 2$, C stands for the absolute constant.

Let $\Omega^a$ denote the set of patterns, located in square axa. Then from the definition of measure $\lambda$ and the above given results, it follows:

<u>Corollary.</u> Measure of the set $\Omega^a$ is finite:

$$\lambda(\Omega^a) = \sum_{n=0}^{\infty} a^n \lambda(\Omega_n) < \infty .$$

58

To prove our theorem we divide sets $\mathcal{Q}_n$ ($n \in \mathbb{N}$) into two subsets

$$\mathcal{Q}_n = \mathcal{Q}_n' \cup \mathcal{Q}_n''$$

where $\mathcal{Q}_n'$ consists of patterns, the total sum $l_n$ of lenghts of the finite straight lines satisfies $l_n < c_n$, $\mathcal{Q}_n''$ consists of patterns where $l_n \geqslant c_n$ (value of $c_n$ is given below). Then

$$\lambda(\mathcal{Q}_n) = \lambda(\mathcal{Q}_n') + \lambda(\mathcal{Q}_n'') \ .$$

To prove the theorem, we shall use special coordinates, which give us a new measure $\mu$ . Later on we shall demonstrate that $\lambda \leqslant \mu$, so the estimate for $\mu(\mathcal{Q}_n)$ holds true for $\lambda(\mathcal{Q}_n)$ also.

We shall evaluate $\lambda(\mathcal{Q}_n')$ and $\lambda(\mathcal{Q}_n'')$ separately.

## 2. Evaluation of the measure of $\mathcal{Q}_n'$

Let us assume that $l_n < c_n$ . We shall call the left end of each finite straight line of the pattern the beginning and the right end the end, and fix a positive integer k. We divide the square into $2^k$ equal ribbons by vertical lines and supply the ribbons with numbers $1, 2, \ldots, 2^k$ from the left to the right.

Let $\mathcal{Q}_n^k \subset \mathcal{Q}_n'$ be a set of patterns, which satisfy the following conditions:

a) none of the finite straight lines has its beginning and end in the same ribbon;

b) if a finite straight line has its beginning in a fixed ribbon, there is not any other finite straight line which has its beginning or end on that finite straight line in this fixed ribbon;

c) no finite straight line begins from the lower and upper side of the first ribbon.

Obviously $\quad \mathcal{Q}_n^1 \subset \mathcal{Q}_n^2 \subset \ldots \subset \mathcal{Q}_n^k \subset \ldots \subset \mathcal{Q}_n'$
$$\mathcal{Q}_n' = \bigcup_{k=1}^{\infty} \mathcal{Q}_n^k \qquad\qquad (n \in \mathbb{N}) \qquad\qquad ( 1 )$$

Let the coordinates of the finite straight lines which

59

begin from the left side of the first ribbon, be $(\theta, p)$, where $\theta$ is the polar angle described above, $p$ is the distance of the beginning point of the finite straight line from point 0 (fig.3).
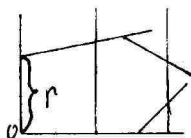


Fig.3.

We arrange the finite straight lines in consecutive order by the growth of coordinate $p$. We determine the coordinate $p$ for a finite straight line which begins in the ribbon numbered $i$ as

$$p = s_{i-1} + l_{i-1} + l_i' + x ,$$

where

$s_{i-1} = 1 + 2(i-1)2^{-k}$, i.e. the perimeter of the part of the square which is covered by the first $i-1$ ribbons minus 1;

$l_{i-1}$ is the summary length of the finite straight lines and the parts of the finite straight lines which are located in the first $i-1$ ribbons;

$l_i'$ equals zero, if the finite straight line under observation begins from the lower side of the ribbon. Otherwise $l_i'$ is the summary length of the parts of the finite straight lines below the finite straight line from which begins the finite straight line under observation, plus $2^{-k}$;

$x$ is the distance of the beginning point of the finite straight line from the left side of the $i$th ribbon, measured along the line from which it begins (fig.4a); for the finite straight lines which begin from the upper or lower side of the ribbon $x$ is measured along that side (fig.4b).
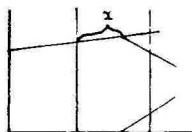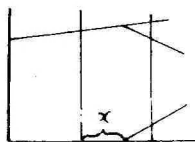


Fig.4a.



Fig.4b.

It is obvious that the coordinates defined as above satisfy the inequalities

$$0 < p_1 < p_2 < \ldots < p_n < 3 + l_n.$$

Thus we have determined the beginning points and the directions of the finite straight lines. To describe the pattern uniquely we must add complementary discreet coordinates $\delta_j$ to show in every junction point of two finite straight lines which one of them ends. We arrange the junction points in consecutive order moving along the ribbons from the bottom to the top and moving from ribbon to ribbon from the left to the right and supply the junction point with coordinate 0 if the lower finite straight line ends, and 1 if otherwise.

Now our pattern is described by the coordinates

$$\left\{ (\theta_1, p_1), \ldots, (\theta_n, p_n); \delta_1, \ldots, \delta_s \right\}, \quad \text{where}$$

$$\left. \begin{array}{l} -\dfrac{\pi}{2} \leqslant \theta_i < \pi \quad (i=1,\ldots,n) \\[2mm] 0 \leqslant p_1 < \ldots < p_n < 3 + l_n \end{array} \right\} \qquad (2)$$

$\delta_j$ equals 0 or 1, $1 \leqslant j \leqslant s$, $s \leqslant n$.

Defining

$$d\mu = d\theta_1 dp_1 \ldots d\theta_n dp_n,$$

we obtain an upper bound for the measure of $\Omega_n^k$ :

$$\mu(\Omega_n^k) = \sum_{\delta_1, \ldots, \delta_1 = 0}^{1} \int \ldots \int d\theta dp <$$

$$< 2^n \left(\frac{3\pi}{2}\right)^n \frac{(3+l_n)^n}{n!} = (3\pi)^n \frac{(3+l_n)^n}{n!},$$

where the integration area is determined by (2).

The estimate obtained above is independent of $k$ and because of (1) holds true for $\Omega_n$ also:

$$\mu(\Omega_n) < (3\pi)^n \frac{(3+l_n)^n}{n!} \qquad (3)$$

Now we shall demonstrate that $\lambda \leqslant \mu$. The transition Jacobian

61

between coordinates $(\theta_1,\ldots,\theta_n,\ \varsigma_1,\ldots,\ \varsigma_n)$ and $(\theta_1,\ldots,\theta_n,p_1,\ldots,p_n)$ is

$$J = \det \begin{bmatrix} \left(\dfrac{\partial\theta_i}{\partial\theta_j}\right)^n_{i,j=1} & \left(\dfrac{\partial\theta_i}{\partial p_j}\right)^n_{i,j=1} \\[2ex] \left(\dfrac{\partial\varsigma_i}{\partial\theta_j}\right)^n_{i,j=1} & \left(\dfrac{\partial\varsigma_i}{\partial\varsigma_j}\right)^n_{i,j=1} \end{bmatrix} = $$

$$ = \det \begin{bmatrix} I & O \\[2ex] \left(\dfrac{\partial\varsigma_i}{\partial\theta_j}\right)^n_{i,j=1} & \triangle \end{bmatrix} = \prod_{i=1}^{n} \dfrac{\partial\varsigma_i}{\partial p_i} \qquad (4)$$

Here I stands for n×n unit matrix (coordinates $\theta_i$ and $\theta_j$ (i≠j) are independent) and O stands for n×n zero-matrix (coordinates $\theta_i$ are independent of coordinates $p_j$). The n×n matrix $\triangle$ is lower-triangular matrix with diagonal elements $\partial\varsigma_i/\partial p_i$ ($\varsigma$-coordinates of the finite straight lines are independent of the p-coordinates of the finite straight lines with bigger numbers).

Let us observe the increment $\partial p$ of the coordinate p of a finite straight line (0 is fixed). Accordingly the increment of the coordinate $\varsigma$ is $\partial\varsigma$ (fig.5).

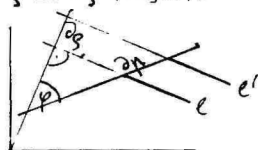

Fig.5.

As we can see,

$$\frac{\partial\varsigma}{\partial p} = \cos\varphi, \text{ so } \frac{\partial\varsigma_i}{\partial p_i} \leqslant 1 \ (i=1,\ldots,n),$$

and after (4)

$$J = \prod_{i=1}^{n} \frac{\partial\varsigma_i}{\partial p_i} \leqslant 1,$$

which proves that $\lambda \leqslant \mu$. So the estimates obtained for the measure $\mu$ hold true for the measure $\lambda$ as well.

Now let us estimate the measure $\Lambda(\Omega'_n)$ proceeding from inequality (3) and give exact values to so far undefined constants $c_n$.

1) Assuming that $l_n \leqslant 1$

$$\Lambda(\Omega'_n) \leqslant \mu(\Omega'_n) < (3\pi)^n \frac{(3+l_n)^n}{n!} \leqslant (12\pi)^n \frac{1}{n!} = \frac{C^n}{n!} \quad , \quad (5)$$

as $3+l_n \leqslant 4$. Here and below $C$ stands for absolute constants.

2) Assuming that $1 < l_n < c_n$, where

$$c_n = n^\alpha \qquad\qquad (6)$$

(value of $\alpha = \alpha(n)$ will be given later), we have $3+l_n < 4l_n < 4n^\alpha$ and according to the inequalities obtained from the Stirling formula

$$C_1 \frac{n^{n+\frac{1}{2}}}{e^n} < n! < C_2 \frac{n^{n+\frac{1}{2}}}{e^n} \quad ,$$

we come to the conclusion that

$$\Lambda(\Omega'_n) \leqslant \mu(\Omega'_n) \leqslant (3\pi)^n \frac{(3+l_n)^n}{n!} < (3\pi)^n \frac{(4n^\alpha)^n}{n!} =$$

$$= (12\pi)^n \frac{n^{n\alpha}}{n!} < (12\pi)^n \frac{e^n n^{n\alpha}}{n^{n+1/2}} < C^n n^{-(1-\alpha)n}$$

Taking

$$\alpha = \frac{\ln(3n \ln\ln n) - \ln\ln n}{\ln n} \qquad , \quad (7)$$

we have

$$C^n n^{-(1-\alpha)n} = C^n \left(\frac{\ln \ln n}{\ln n}\right)^n .$$

So, if $c_n$ is determined by (6) and (7), then

$$\Lambda(\Omega'_n) < C^n \left(\frac{\ln \ln n}{\ln n}\right)^n . \qquad\qquad (8)$$

## 3. Evaluation of the measure of $\Omega''_n$

Here we consider patterns which belong to $\Omega''_n$, i.e. total sum of lengths of the finite straight lines of which $l_n \geqslant c_n = n^\alpha$ ($\alpha$ is determined by (7)). To begin with let us supply the finite straight lines with the same coordinates

as in p.2.

Now let us observe the finite straight lines which are longer than

$$\frac{1}{2}\frac{n^{\alpha}}{n} = \frac{1}{2}n^{\alpha-1}$$

Obviously, the number of these lines m must satisfy the inequality

$$\frac{n^{\alpha}}{2n}(n-m) + \sqrt{2}m > n^{\alpha},$$

(as $\sqrt{2}$ is the maximum possible length of a finite straight line in unit square). Thus we obtain

$$m > \frac{n}{2\sqrt{2}n^{1-\alpha} - 1}$$

From now on we assume that at least half of these finite straight lines have

$$\theta \in [\frac{\pi}{4}, \frac{3\pi}{4}].$$

If that does not hold true, then by turning our square by $\frac{\pi}{2}$ we can achieve it, nevertheless. So the measure of the set of the patterns which satisfy this condition cannot be more than twice smaller than the measure of $\Omega_n''$. Let $H \subset \{1, \ldots, n\}$ be the set of indexes of these finite straight lines the lenght of which is more than $\frac{1}{2}n^{\alpha-1}$ and

$$\theta \in [\frac{\pi}{4}, \frac{3\pi}{4}].$$

Obviously, the number of these finite straight lines satisfies

$$|H| \geqslant \left[\frac{m}{2}\right] \geqslant \left[\frac{n}{2\sqrt{2}n^{1-\alpha} - 1}\right] = p, \qquad (9)$$

Let us supply p lines with new coordinates. For that we divide the square using vertical lines into ribbons with the width

$$\frac{n^{\alpha-1}}{4\sqrt{2}}$$

Every finite straight line of those with $\theta \in [\frac{\pi}{4}, \frac{3\pi}{4}]$ intersects both sides of one ribbon at least. Assembling the ribbons end to end we obtain one ribbon with the length $L = 4\sqrt{2}n^{1-\alpha}$. Let the coordinates of the finite straight lines be:

64

$q_i$ - the distance of the left intersection point from the lower end of the ribbon;

$r_i$ - the distance of the right intersection point from the lower end of the ribbon.

Define $\mu$ by

$$d\mu = dq_1 dr_1 \ldots dq_p dr_p.$$

The following inequalities hold true:

$$0 \leqslant q_1 < q_2 < \ldots < q_p \leqslant L$$
$$0 \leqslant r_1 < r_2 < \ldots < r_p \leqslant L \, ,$$

and the 2p-dimensional $\mu$-measure of this set of finite straight lines equals

$$\int \ldots \int_{A'} \prod_{i \in H} dq_i dr_i = \frac{L^{2p}}{(p!)^2}$$

Coordinates $(q_i, r_i)$ and polar coordinates $(\theta_i, \zeta_i)$ are related by

$$\theta_i = \frac{\pi}{2} + \arctan L(r_i - q_i) \, ,$$

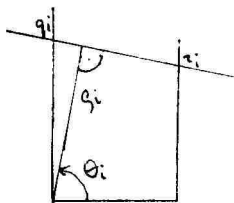$$\zeta_i = \frac{r_i}{\sqrt{1 + L^2(r_i - q_i)^2}} \, ,$$



Fig.6.

and the transition Jacobian can be evaluated by

$$J = \det \begin{bmatrix} \dfrac{\partial \zeta_i}{\partial r_i} & \dfrac{\partial \zeta_i}{\partial q_i} \\[2mm] \dfrac{\partial \theta_i}{\partial r_i} & \dfrac{\partial \theta_i}{\partial q_i} \end{bmatrix} = \frac{L}{[1 + L^2(r_i - q_i)^2]^{3/2}} \leqslant L \, .$$

So the 2p-dimensional $\lambda$-measure of the set of these finite straight lines satisfies the inequality

$$\int \ldots \int_{A'} \prod_{i \in H} d\zeta_i d\theta_i \leqslant L^p \mu(A') = \frac{L^{3p}}{(p!)^2}$$

If we fix set H and coordinates $\zeta_i, \theta_i$ ($i \in H$), then for the coordinates of the rest of the finite straight lines (2)

65

holds true. So the $2(n-p)$ dimensional measure of the corresponding intersection of the set $\Omega_n$ is not more than

$$2^n \left(\frac{3\pi}{2}\right)^{n-p} \frac{(3+1_n)^{n-p}}{(n-p)!} \ ,$$

and since there are $C_n^p$ possibilities for locating a set $H$, we have

$$\Lambda(\Omega_n^{\shortparallel}) < 2 C_n^p \ 2^n \ \left(\frac{3\pi}{2}\right)^{n-p} \frac{(3+1_n) \ L^{3p}}{(n-p)! \ (p!)^2} \ =$$

$$= 2^{p+1} (3\pi)^{n-p} L^{3p} \frac{(3+1_n)^{n-p} n!}{[(n-p)!]^2 \ (p!)^3} \ <$$

(since $3+1_n < 2+\sqrt{2} n < 5n$)

$$< c^n \ n^{3(1-\alpha)p} \ \frac{(5n)^{n-p} \ n!}{[(n-p)!]^2 (p!)^3} \ <$$

(using Stirling formula)

$$< c^n \ n^{3(1-\alpha)p} \ \frac{n^{n-p} e^{2(n-p)} e^{-n} n^{n+\frac{1}{2}} e^{-3p}}{(n-p)^{2(n-p)+1} p^{3(p+\frac{1}{2})}}$$

(here we evaluate all the constants in various powers in terms of $c^n$ and substitute the value of $p$ from (9))

$$< c^n \ n^{\frac{2-3\alpha}{2\sqrt{2}} n^\alpha} < c^n \ n^{-\frac{1}{3} n^\alpha} \ ,$$

as

$$\frac{2-3\alpha}{2\sqrt{2}} = \frac{3(1-\alpha)}{2\sqrt{2}} - \frac{1}{2\sqrt{2}} \qquad , \ \alpha = \alpha(n) \uparrow 1 ,$$

we have in case of large values of $n$

$$\frac{3(1-\alpha)}{2\sqrt{2}} - \frac{1}{2\sqrt{2}} < -\frac{1}{3} \ .$$

At last, substituting the value of $\alpha$ from (7), we obtain the estimation

$$\Lambda(\Omega_n^{\shortparallel}) < c^n \left(\frac{1}{\ln n}\right)^n \qquad . \tag{10}$$

Estimates (8) and (10) together prove the theorem.

References

1. Arak,T.,Surgailis,D.    Markov    fields   with   polygonal
   realizations. Probability and Related Fields. To appear
2. Kendall,M.G.,Moran,P.A.P.    Geometerical    probability.
   Griffin's  Statistical   Monographs  and  Courses.   No.5,
   London,1963. C.Griffin

## ОЦЕНКА МЕРЫ МНОЖЕСТВА ПОЛИГОНАЛЬНЫХ МОЗАИКОВ
### П.Микков

Р е з ю м е

   В  данной работе рассматривается множество $\Omega$ мозаиков
на единичном квадрате, состоящих из прямых отрезков и содер-
жащих  узлы  только  типа $\top$. Множество  мозаиков, состоящих
из   $n$  ($n \in \mathbb{N}$) отрезков, обозначается   через $\Omega_n$. На множестве
$\Omega_n$ определяется  мера  $\lambda$ , соответствующая изучаемой в [1]
марковской мере  при нулевом потенциале. Основным результатом
статьи является соотношение
$$\lambda(\Omega) = \sum_{n=0}^{\infty} \lambda(\Omega_n) < \infty,$$
которое  получается  как следствие из следующей оценки меры
множества $\Omega_n$:
$$\lambda(\Omega_n) < c^n \left( \frac{\ln \ln n}{n} \right)^n,$$
где  $c$  абсолютная константа.

# AN ESTIMATION METHOD FOR COEFFICIENTS OF POLYNOMIAL TREND IN TIME SERIES

## V.Joala

### 1. Introduction

A model often used to describe the trend in time series is the following

$$x(t) = a_0 + \sum_{i=1}^{n} a_i t^i + u_t \qquad , \qquad t = 1, \ldots, M . \qquad (1)$$

Here we assume that the unmeasurable $u_t$ are independently and identically distributed uncorrelated random variables with zero mean and variance $\delta^2$ ( $E u_t = 0$ ; $E u_t^2 = \delta^2$ ; $E u_t u_s = 0$, $t \neq s$ ). The constant $n$ is the order of the polynomial trend.

Usually to estimate coefficients $a_i$, $i = 0, \ldots, n$ the least-square method is used, which is the best unbiased linear method in that case. But there are practical tasks with long series (in engineering, physics, scientific experiments, etc.) for which the least-square method works too slowly. It is slow because we must compute sums

$$\sum_{t=1}^{M} x(t) \cdot t^i \quad , \quad i = 0, \ldots, n .$$

To compute them we do $M n(n+1)/2$ multiplications and $M(n+1)$ addition operations. The method which uses only M addition operations instead of that is given below.

### 2. The basic ideas

Let the model (1) hold. We shall partiate the sequence of $x(t), t = 1, \ldots, M$ to the $p$ subsequences with the same length $\kappa = M/p$ (assume $M = \kappa \cdot p$ ). In all subsequences we compute sums

$$S_j = \sum_{t=(j-1)\kappa+1}^{jk} x(t) \quad , \quad j = 1, \ldots, p . \qquad (2)$$

We shall construct functions $f_i$ $(i = 0, \ldots, n)$ which make it possible to estimate coefficients $a_i$ by sums $S_j$

$$\hat{a}_i = f_i(S_1, \ldots, S_p) \quad , \quad i = 0, \ldots, n .$$

The experience with the one-dimensional case suggests that we should use functions in form

$$\hat{a}_i = \frac{1}{R_i} \sum_{j=1}^{\mu} g_{ij} s_j - \psi_i (\hat{a}_{i+1}, \ldots, \hat{a}_n), \qquad (3)$$

where

$R_i$ — the functions of coefficients $g_{ij}$ ;

$g_{ij}$ — coefficients with values $-1, 0, 1$ ;

$\mu$ — a number of sums;

$\psi_i$ — function which eliminated the influence of coefficients $a_{i+1}$ and higher (by $i$ ).

Estimating coefficients $a_i$ by formula (3) we must add elements $x(t), t = 1, \ldots, M$ only once.

## 3. Walsh functions

Sums in formula (3) can be represented as Walsh functions, see[2]. In Fig. 1 graphical forms of coefficients $g_{ij}$ ror $n = 0, \ldots, 3$ are given.
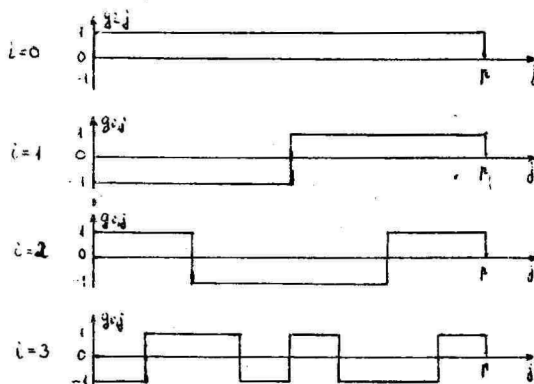


Figure 1. Graphical forms of Walsh functions

For Walsh function we shall use the notation $W(i, 2^k)$. For every $n$ there exists a family of Walsh functions $W(0, 2^k), W(1, 2^k), \ldots, W(n, 2^k)$ . To define a function $W(i, 2^k)$ let us suppose that we have a sequence $\{z_i\}$ with $k \cdot 2^k$ elements. We partiate that sequence to the $2^k$ subsequences with the equal length $k$ and compute sums $s_j$ $(j = 1, \ldots, 2^k)$ by formula (2). Now we show how coefficients $g_{ij}, j = 1, \ldots, 2^k$ are found out.

In $W(0, 2^k)$ all $g_{0,j}=1, j=1, \ldots, 2^k$ . To define $W(1, 2^k)$ from $W(0, 2^k)$ the first half of $g_{0,j}$ (here $j = 1, \ldots, 2^{k-1}$ ) will be inverted and the second half will not change (See Fig. 1).

To define $W(2,2^n)$ from $W(1,2^n)$ the first and the third quarters of $q_{i,j}$ (by $j$) from $W(1,2^n)$ will be inverted, etc. Some Walsh functions $W(i,\mu)$ for fixed $n$ ($\mu = 2^n$) are given below ($n = 0, \ldots, 3$).

$n = 0$     $W(0,1) = S_1$

$n = 1$     $W(0,2) = S_2 + S_1$
             $W(1,2) = S_2 - S_1$

$n = 2$     $W(0,4) = S_4 + S_3 + S_2 + S_1$
             $W(1,4) = S_4 + S_3 - S_2 - S_1$
             $W(2,4) = S_4 - S_3 - S_2 + S_1$

$n = 3$     $W(0,8) = S_8 + S_7 + S_6 + S_5 + S_4 + S_3 + S_2 + S_1$
             $W(1,8) = S_8 + S_7 + S_6 + S_5 - S_4 - S_3 - S_2 - S_1$
             $W(2,8) = S_8 + S_7 - S_6 - S_5 - S_4 - S_3 + S_2 + S_1$
             $W(3,8) = S_8 - S_7 - S_6 + S_5 - S_4 + S_3 + S_2 - S_1$

If we have sequence $\{t^m\}$ where $t = 1, \ldots, K \, 2^n$ and we compute sums $S_j$ by formula (2), then

$$W(i, 2^n) = 0 \quad , \quad i < m. \qquad (4)$$

Coefficients $a_i$ are estimated by formula

$$\hat{a}_i = \frac{W(i, 2^n)}{R_n} - f_i(\hat{a}_{i+1}, \ldots, \hat{a}_n) \quad , \quad i = 0, \ldots, n \quad ,$$

where

$$R_i = \frac{M^{n+1}}{(n+1)(\prod_{j=0}^{n} 2^j)^2} \sum_{j=1}^{2^n} g_{ij} \left[ j^{n+1} - (j-1)^{n+1} \right] . \qquad (5)$$

Using property (4) it is easy to see, that we can choose the functions $f_i$ so, that $E\hat{a}_i = a_i$. For coefficient $\hat{a}_n$ we can see that

$$D(\hat{a}_n) = \frac{M}{(R_n)^2} \delta^2 .$$

To attain subsequences with an equal length we cannot use some elements of $x(t)$ or choose $M = K \cdot 2^n$, $K = 1, 2, \ldots$ in long series.

### 4. Modified functions

In the paper [1] it is shown that variances $D(\hat{a}_n)$ will be minimized by modifying functions $W(i, \mu)$ . Therefore we must compute more sums $S_j$ ($j = 1, \ldots, 1.5 \mu$) than in the

case of Walsh functions, but we win in accuracy of estimating and do not lose in computing time.

In Fig. 2 graphical forms of coefficients modified
functions for $i = 0, \ldots, 3$ (here $\rho^\varkappa = 1{,}5 \cdot \rho$ ) are given.
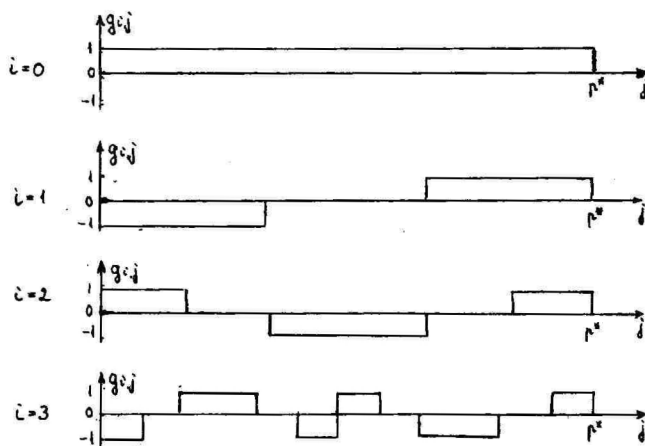


Figure 2. Graphical forms of modified functions

In Figure 2. it is shown that for all $i$ (except $i = 0$ )
to estimate $a_i$ one third of sums $S_j$ is not used.

To define modified functions $V(i, \rho^\varkappa)$ we must use the
same algorithm which was used for Walsh functions (here $\rho^\varkappa =$
$= 1{,}5 \cdot 2^n$ ) up to $i = n-1$ . To define $V(n, \rho^\varkappa)$ from $V(n-1, \rho^\varkappa)$
we must partiate all subsequences in $V(n-1, \rho^\varkappa)$ into three
parts, not two. The first of them will be inverted the second will not be used (here $g_{i,j} = 0$ ) and the third of them
will not be changed.

Modified functions $V(i, \rho^\varkappa)$ for estimating coefficients
$a_i$ are given below ( $n = 0, \ldots, 3$ ).

$$n = 0 \qquad V(0, 1) = S_1$$

$$n = 1 \qquad V(0, 3) = S_3 + S_2 + S_1$$

$$V(1, 3) = S_3 - S_1$$

71

$$n=2 \qquad V(0,6) = S_6 + S_5 + S_4 + S_3 + S_2 + S_1$$
$$V(1,6) = S_6 + S_5 - S_2 - S_1$$
$$V(2,6) = S_6 - S_4 - S_3 + S_1$$

$$n=3 \qquad V(0,12) = S_{12} + S_{11} + S_{10} + S_9 + S_8 + S_7 + S_6 + S_5 + S_4 + S_3 + S_2 + S_1$$
$$V(1,12) = S_{12} + S_{11} + S_{10} + S_9 - S_4 - S_3 - S_2 - S_1$$
$$V(2,12) = S_{12} + S_{11} - S_8 - S_7 - S_6 - S_5 + S_2 + S_1$$
$$V(3,12) = S_{12} - S_{10} - S_9 + S_7 - S_6 + S_4 + S_3 - S_1$$

It is easy to prove that property (4) holds for all $V(i, p^x)$ too. Coefficients $a_i$ are estimated by formula

$$\bar{a}_i = \frac{V(i, 1.5 \cdot 2^n)}{P_i} - \varphi_i(\bar{a}_{i+1}, \dots, \bar{a}_n) \quad , \quad i = 0, \dots, n \, ,$$

where $P_i = 1,125 \cdot R_i$ and $R_i$ are calculated by formula (5).

Using property (4) it is easy to see, that we can choose the functions $\varphi_i$ so, that $E\bar{a}_i = a_i$. For coefficients $\bar{a}_n$ we can see that

$$\mathcal{D}(\bar{a}_n) = \frac{2M}{3(P_n)^2} \mathcal{b}^2.$$

When $n=1$, coefficients $\bar{a}_1$ and $\bar{a}_0$ have respectively 12,5% and 9,4% greater variances than the coefficients estimated by the least-square method.

To attain subsequences with an equal length here, one can pass some elements of $x(t)$ or choose $M = 1,5 \cdot K \cdot 2^n$, $K = 1, 2, \dots$ in long series.

## 5. Computation formulae

$$n=0 \qquad x(t) = a_0 + u_t \quad , \quad t = 1, \dots, M$$
$$\bar{a}_0 = \frac{1}{M} \sum_{t=1}^{M} x(t) = \frac{1}{M} S_1$$

$$n=1 \qquad x(t) = a_0 + a_1 t + u_t \quad , \quad t = 1, \dots, M$$
$$\bar{a}_1 = \frac{4,5}{M^2} (S_3 - S_1)$$
$$\bar{a}_0 = \frac{1}{M} (S_3 + S_2 + S_1 - \bar{a}_1 \frac{M(M+1)}{2})$$
$$S_j = \sum_{t=(j-1)K+1}^{jK} x(t) \quad , \quad K = \frac{M}{3} \quad , \quad j = 1, \dots, 3$$

72

$n=2$  $x(t) = a_0 + a_1 t + a_2 t^2 + u_t$ , $t = 1, \ldots, M$

$$\bar{a}_2 = \frac{18}{M^3}(S_6 - S_4 - S_3 + S_1)$$

$$\bar{a}_1 = \frac{4.5}{M^2}\left(S_6 + S_5 - S_2 - S_1 - \bar{a}_2\frac{M^2(M+1)}{4.5}\right)$$

$$\bar{a}_0 = \frac{1}{M}\left(S_6 + S_5 + S_4 + S_3 + S_2 + S_1 - a_2\frac{M(M+1)(2M+1)}{6} - \bar{a}_1\frac{M(M+1)}{2}\right)$$

$$S_j = \sum_{t=(j-1)k+1}^{jk} x(t) \quad , \quad k = \frac{M}{6} \quad , \quad j = 1, \ldots, 6$$

$n=3$  $x(t) = a_0 + a_1 t + a_2 t^2 + a_3 t^3 + u_t$ , $t = 1, \ldots, M$

$$\bar{a}_3 = \frac{96}{M^4}(S_{12} - S_{10} - S_9 + S_7 - S_6 + S_4 + S_3 - S_1)$$

$$\bar{a}_2 = \frac{18}{M^3}\left(S_{12} + S_{11} - S_8 - S_7 - S_6 - S_5 + S_2 + S_1 - \bar{a}_3\frac{M^3(M+1)}{12}\right)$$

$$\bar{a}_1 = \frac{4.5}{M^2}\left(S_{12} + S_{11} + S_{10} + S_9 - S_4 - S_3 - S_2 - S_1 - \right.$$
$$\left. - \bar{a}_2\frac{M^2(M+1)}{4.5} - \bar{a}_3\frac{M^2(16M^2+27M+9)}{81}\right)$$

$$\bar{a}_0 = \frac{1}{M}\left(S_{12} + S_{11} + S_{10} + S_9 + S_8 + S_7 + S_6 + S_5 + S_4 + S_3 + S_2 + S_1 - \right.$$
$$\left. - \bar{a}_1\frac{M(M+1)}{2} - \bar{a}_2\frac{M(M+1)(2M+1)}{6} - \bar{a}_3\frac{M^2(M+1)^2}{4}\right)$$

$$S_j = \sum_{t=(j-1)k+1}^{jk} x(t) \quad , \quad k = \frac{M}{12} \quad , \quad j = 1, \ldots, 12.$$

## 6. Computer experiments

The third-order model was used to examine how this method works. Here
$$x(t) = a_0 + a_1 t + a_2 t^2 + a_3 t^3 + u_t \quad , \quad t = 1, \ldots, M.$$

There were three different trend models (coefficients of them are given in Table 1.)

Table 1.

Simulated trends coefficients

| : | : | I | : | II | : | III | : |
|---|---|---|---|---|---|---|---|
| : $a_0$ | : | 50. | : | -1200. | : | 22.505 | : |
| : $a_1$ | : | -6. | : | 170. | : | 0.728 | : |
| : $a_2$ | : | 0.02 | : | -2. | : | 0.0907 | : |
| : $a_3$ | : | 0.00004 | : | 0.0055 | : | -0.0000739 | : |

73

Random values $u_t$ are generated by a random number generator on Computer POP-11. There were random values $u_t$ with the uniform distribution (variances 80000, 800, 8) and with the normal distribution (variances 7000, 70, 0.7).

The error measure was

$$Z_i = 100 \cdot \frac{E[x_i(t) - \hat{x}_i(t)]^2}{E[x_i(t) - \overline{x}_i(t)]^2} \qquad , \quad i = 1, \ldots, 50 \ ,$$

where $x(t)$ - generated series, $\overline{x}(t)$ - estimated trend by least-square method, $\hat{x}(t)$ - estimated trend by fast method. The length of the series was M=240 and estimation has been made 50 times.

$$MEAN = \frac{1}{50} \sum_{i=1}^{50} Z_i \ ,$$

$$MAX = \max_i Z_i \ , \quad i = 1, \ldots, 50 \ .$$

In the Tables 2.-5. the computer experiment results are given.

Table 2

MEAN: Uniform distribution

| : Variance : | I. : | II | : | III | : |
|---|---|---|---|---|---|
| : 80000 | : 100.53: | 100.59 | : | 100.63 | : |
| : 800 | : 100.64: | 100.67 | : | 100.75 | : |
| : 8 | : 100.55: | 100.67 | : | 100.58 | : |

Table 3

MAX: Uniform distribution

| : Variance : | I | : | II. | : | III | : |
|---|---|---|---|---|---|---|
| : 80000 | : 106.20 | : | 103.18 | : | 103.46 | : |
| : 800 | : 102.76 | : | 103.49 | : | 103.18 | : |
| : 8 | : 103.42 | : | 105.05 | : | 106.20 | : |

Table 4

MEAN: Normal distribution

| : Variance : | I | : | II | : | III | : |
|---|---|---|---|---|---|---|
| : 7000 | : 100.49 | : | 100.44 | : | 100.46 | : |
| : 70 | : 100.58 | : | 100.48 | : | 100.65 | : |
| : 0.7 | : 100.50 | : | 100.67 | : | 100.51 | : |

MAX: Normal distribution

| : Variance : | I | : | II | : | III | : |
|---|---|---|---|---|---|---|
| : 7000 | : 101.89 | : | 101.32 | : | 101.96 | : |
| : 70 | : 102.19 | : | 103.22 | : | 102.38 | : |
| : 0.7 | : 102.08 | : | 102.03 | : | 103.52 | : |

## 7. Conclusions

In this paper an estimation method of coefficients of polynomial trend in time series is given. It works much faster on computers than the least-square method and experimentally it is shown that in case of accuracy it is approximately the same as the least-square method. Therefore this method is useful to increase the productivity of computer programs, especially in case of long time series or automatic control systems where results must be computed very fast.

References

1. Joala V., Olman V. An Estimation Method of Coefficients of One-dimensional Linear Regression. Proceedings of the Academy of Siences of the Estonian SSR, Physics Mathematics, 1987, №4, 422-424

2. Проектирование специализированных информационно-вычислительных систем (Под ред. Ю.Смирнова). М., 1984

МЕТОД ОЦЕНИВАНИЯ КОЭФФИЦИЕНТОВ ПОЛИНОМИАЛЬНОГО
ТРЕНДА ВРЕМЕННЫХ РЯДОВ
В.Йоала

Р е з ю м е

В данной статье приведен метод оценивания коэффициентов полиномиального тренда временных рядов, который на вычислительных машинах работает быстрее, чем метод наименьших квадратов. Экспериментально показано, что по точности предлогаемый метод сравним с методом наименьших квадратов. Этот метод можно применять для повышения производительности программ анализа временных рядов на ЭВМ, особенно, когда временные ряды длинные, а также в системах автоматического управления, где результаты должны быть вычислены максимально быстро.

# CONTENTS