ISSN 1407-5806

# COMPUTER DODELLING AND NEW TECHNOLOGIES

Volume 14 No 4

2010

# Computer Modelling and New Technologies

Volume 14, No.4 – 2010

**ISSN 1407-5806 ISSN 1407-5814** (On-line: www.tsi.lv)

**Riga – 2010** 

#### **EDITORIAL BOARD:**

Prof. Igor Kabashkin (Chairman of the Board), *Transport & Telecommunication Institute, Latvia;*Prof. Yuri Shunin (Editor-in-Chief), *Information Systems Management Institute, Latvia;*Prof. Adolfas Baublys, *Vilnius Gediminas Technical University, Lithuania;*Dr. Brent D. Bowen, *Purdue University, USA;*Prof. Olgierd Dumbrajs, *Helsinki University of Technology, Finland;*Prof. Eugene Kopytov, *Transport & Telecommunication Institute, Latvia;*Prof. Arnold Kiv, *Ben-Gurion University of the Negev, Israel;*Prof. Juris Zakis, *University of Latvia;*Prof. Edmundas Zavadskas, *Vilnius Gediminas Technical University, Lithuania.*

#### **Host Organization:**

Transport and Telecommunication Institute

#### Supporting Organizations:

Latvian Transport Development and Education Association Latvian Academy of Sciences Latvian Operations Research Society

#### THE JOURNAL IS DESIGNED FOR PUBLISHING PAPERS CONCERNING THE FOLLOWING FIELDS OF RESEARCH:

- mathematical and computer modelling
- mathematical methods in natural and engineering sciences
- physical and technical sciences
- computer sciences and technologies
- semiconductor electronics and semiconductor technologies
- aviation and aerospace technologies
- electronics and telecommunication
- navigation and radar systems
- telematics and information technologies
- transport and logistics
- economics and management
- social sciences In journal articles can be presented in English. All articles are reviewed.

#### EDITORIAL CORRESPONDENCE

Transporta un sakaru institūts (Transport and Telecommunication Institute) Lomonosova iela 1, LV-1019, Riga, Latvia. Phone: (+371) 67100593. Fax: (+371) 67100535. E-mail: journal@tsi.lv, www.tsi.lv

#### COMPUTER MODELLING AND NEW TECHNOLOGIES, 2010, Vol. 14, No.4

**ISSN** 1407-5806, **ISSN** 1407-5814 (on-line: www.tsi.lv) Scientific and research journal of Transport and Telecommunication Institute (Riga, Latvia) The journal is being published since 1996.

Copyright <sup>©</sup> Transport and Telecommunication Institute, 2010

# CONTENTS

Editors' remarks	5					
Stochastic Models in Reliability Engineering, Life Sciences and Operations Management						
Z. Laslo, G. Gurevich	6					
Applied Statistics and Operations Research	7					
Residual Partial Sums Techniques for Fixed Designs to Find Change-Points in Linear Regression W. Bischoff	7					
Beta-Distribution Models in Stochastic Project Management D. Greenberg, A. Ben-Yair	14					
Continuous Models of Current Stock of Divisible Productions E. Kopytov, S. Guseynov, E. Puzinkevich, L. Greenglaz	19					
Statistical Inference Using Entropy Based Empirical Likelihood Statistics G. Gurevich, A. Vexler	31					
for Interval Censored and Truncated Data F. Vonta, C. Huber	40					
Upon Controlling Several Building Projects in a Two-Level Construction System D. Golenko-Ginzburg, Z. Laslo	50					
Performance Evaluation of Teamwork Reflexive Analysis G. I. Petkov, V. A. Iliev	57					
<b>The Hotelling's Metric as a Cluster Stability Measure</b> Z. Volkovich, Z. Barzily, D. Toledano-Kitai, R. Avros	65					
Predictions on the Formation of Nerve-Muscle Connection I. Nowik	73					
Authors' Index	77					
Personalia	78					
Cumulative Index	80					
Preparation of Publications						



### Editors' Remarks

### The Paradox of Time

*By Henry Austin Dobson* 

Time goes, you say? Ah no! Alas, Time stays, we go; Or else, were this not so, What need to chain the hours, For Youth were always ours? Time goes, you say?- ah no!

Ours is the eyes' deceit Of men whose flying feet Lead through some landscape low; We pass, and think we see The earth's fixed surface flee:-Alas, Time stays,- we go!

Once in the days of old, Your locks were curling gold, And mine had shamed the crow. Now, in the self-same stage, We've reached the silver age; Time goes, you say?- ah no! Once, when my voice was strong, I filled the woods with song To praise your "rose" and "snow"; My bird, that sang, is dead; Where are your roses fled? Alas, Time stays,- we go!

See, in what traversed ways, What backward Fate delays The hopes we used to know; Where are our old desires?-Ah, where those vanished fires? Time goes, you say?- ah no!

How far, how far, O Sweet, The past behind our feet Lies in the even-glow! Now, on the forward way, Let us fold hands, and pray; Alas, Time stays,- we go!

\*\*\*\*\*\*\*\*\*\*

# Henry Austin Dobson (1840-1921)

This 14<sup>th</sup> volume No.4 presents the special selection of articles of the International Symposium on Stochastic Models in Reliability Engineering, Life Sciences and Operations Management [SMRLO'10].

Our journal policy is directed on the fundamental and applied sciences researches, which are the basement of a full-scale modelling in practice.

This edition is the continuation of our publishing activities. We hope our journal will be interesting for research community, and we are open for collaboration both in research and publishing. This number continues the current 2010 year of our publishing work. We hope that journal's contributors will consider the collaboration with the Editorial Board as useful and constructive.

EDITORS

In Summin\_

Yu.N. Shunin

I.V. Kabashkin

#### GUEST EDITORS' REMARKS

This Special Issue is devoted to papers from the International Symposium on Stochastic Models in Reliability Engineering, Life Sciences and Operations Management [SMRLO'10]. The Symposium was a continuation of the International Symposium on Stochastic Models in Reliability, Safety, Security and Logistics (SMRSSL'05) held in 2005 and also was held at the SCE-Shamoon College of Engineering, Beer Sheva, Israel on February 2010.

The idea of the Symposium was to assemble researchers and practitioners from universities, institutions, industries, businesses and government, working in these fields. Theoretical issues and applied case-studies, presented on the symposium, were ranged from academic considerations to operational applications.

Presenters come from more than thirty different countries all around the world: Bulgaria, Canada, China, Czech Republic, Cyprus, France, Germany, Greece, India, Ireland, Israel, Italia, Latvia, Lithuania, Netherlands, Norway, U.K., Poland, Portugal, Romania, Russia, Serbia, Singapore, Slovakia, South Africa, Spain, Sweden, Switzerland, Taiwan, Turkey, Ukraine, and USA.

One hundred and eighty five papers were accepted for presentation at the conference and publication in the Symposium Proceedings. Later on these articles were reviewed for possible extension and inclusion in the journal. Authors of nine of the articles were invited to submit their studies for publication in this Special Issue of **Computer Modelling and New Technologies**.

We hope that this selection of papers will give an idea of the diversity of topics covered in the SMRLO'10.

Prof. Zohar Laslo and Dr. Gregory Gurevich

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 7–13 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# RESIDUAL PARTIAL SUMS TECHNIQUES FOR FIXED DESIGNS TO FIND CHANGE-POINTS IN LINEAR REGRESSION

#### W. Bischoff

Catholic University Eichstätt-Ingolstadt, Faculty of Mathematics and Geography Ostenstr. 26–28, 85072 Eichstätt, Germany E-mail: wolfgang.bischoff@ku-eichstaett.de

We investigate a data set describing the quality of a production process. By the information of these data it has to be decided whether the quality is constant or whether the quality changes. Our null hypothesis is that the quality is constant that is a linear regression. In practice it is popular to investigate the partial sums of the least squares residuals to look for changes in linear regression. The partial sums of the least squares residuals can be embedded into the class of continuous functions. By this procedure we obtain a stochastic process with continuous paths. It is called residual partial sum process. If the number of observations is large enough a projection of the Brownian motion can be considered as approximation (with respect to weak convergence) of the residual partial sum process. This projection of the Brownian motion can be used to establish non-parametric tests of Cramér-von Mises and Kolmogorov–Smirnov type to test for changes in linear regression. We use this procedure to test the data for constant quality.

**Keywords:** residual partial sum limit processes, linear regression models, fixed designs, Brownian motion, projections of Brownian motion, reproducing kernel Hilbert space, change-point problem

#### 1. Introduction

We investigate a data set from a company that provides car companies with toothed lock washers. The company has to guarantee the quality of each delivery of contract goods. For that at the end of the production process a sample of the goods is taken at equidistant time points to check the quality. Knowing the corresponding data the company has to decide whether the quality is constant or whether there is a change. The data of the company are shown on Figure 1 together with the least squares estimation of a constant function.



*Figure 1.* Data to be checked for constant quality. The constant line is the least squares estimation of a constant quality

More general we consider the problem whether a change of a regression occurs during the period of time (or in the experimental region) [0,1], say. To explain the problem in more detail, let the true regression model be given by  $Y(t) = g(t) + \mathcal{E}(t)$ ,  $t \in [0,1]$ , where g is the true but unknown regression function,  $\mathcal{E}(t)$  is a real random variable with expectation 0 and variance  $\sigma^2 > 0$ . We are interested in testing whether the regression function g belongs to a linear regression model, i.e. whether there exist suitable constants  $\beta_1, \ldots, \beta_d \in \mathbb{R}$  such that  $g(t) = \sum_{i=1}^d \beta_i f_i(t), t \in [0,1]$ , where  $f_1, \ldots, f_d : [0,1] \to \mathbb{R}$  are known functions. Hence, we look for a test of the hypothesis  $H_0$ : there is no change, that is:

$$g = \sum_{i=1}^{d} \beta_i f_i = f^T \beta \text{ for some } \beta_1, \dots, \beta_d \in \mathbb{R},$$
(1)

where  $f^T = (f_1, \dots, f_d), \beta = (\beta_1, \dots, \beta_d)^T$ , against the alternative *K*: there is a change, that is there are no  $\beta_1, \dots, \beta_d \in \mathbb{R}$ , such that

$$g = \sum_{i=1}^{d} \beta_i f_i.$$
(2)

For that company it is important to know whether the quality keeps constant or is getting worse (or changing) during a certain period of time. Therefore we are interested in testing whether the mean of the quality of the production is constant, thus the null hypothesis is given by

$$H_0: E(Y(t)) = \beta = \mathbf{1}_{[0,1]}(t)\beta, \quad t \in [0,1], \beta \in \mathbb{R} \text{ unknown constant},$$
(3)

where  $\mathbf{1}_{[0,1]}$  is the function identical 1 on [0,1]. Hence, we have under  $H_0$  the model

$$Y(t) = \beta + \varepsilon(t) = \mathbf{1}_{[0,1]}(t)\beta + \varepsilon(t), \quad t \in [0,1].$$
(4)

In the literature on 'detecting change-points' in regression models, it is common to consider residual partial sums processes or variants of it; see for instance, Gardner [13], Brown, Durbin and Evans [12], Sen and Srivastava [23], Sen [22], Jandhyala and MacNeill [15], Jandhyala and MacNeill [16–18], Watson [26], Bischoff [1], Jandhyala, Zacks and El-Shaarawi [19], Jandhyala and Al-Saleh [14], Bischoff and Miller [10], Xie and MacNeill [27]. We follow the approach of MacNeill [20, 21] and Bischoff [2], who considered a (residual) partial sums process approach. It is worth noting that there is a related approach by empirical processes, see Stute [24] and the papers cited there. In those papers random designs are considered. As opposed to this the papers using partial sums take equidistant designs. Bischoff [1], however, investigates arbitrary fixed designs for the partial sums process approach.

In Section 2 we give some assumptions and discuss preliminaries. Then in Section 3 the theoretical main result is stated, namely the limit of the residual partial sum process when the number of observations goes to infinity. This result can be used to build tests for change-points. It is worth mentioning that for this approach the distribution of the error does not have to be specified. Tests of Kolmogorov(–Smirnov) type based on that limit process are considered in detail in Section 4. Moreover, in Section 4 we apply these tests to the data discussed above.

#### 2. Preliminaries and Assumptions

Let  $n_0$  be the number of observations. We assume that the data are taken at the equidistant design points  $t_{n_01} = \frac{1}{n_0}, t_{n_02} = \frac{2}{n_0}, \dots, t_{n_0n_0} = 1$ . These design points can be embedded in a triangular array of designs points:  $t_{n_1} = \frac{1}{n}, t_{n_2} = \frac{2}{n}, \dots, t_{n_n} = 1, n \in \mathbb{N}$ . Given the true model  $Y(t) = g(t) + \varepsilon(t), t \in [0,1]$ , we have a corresponding array of observations:

$$Y(t_{nj}) = g(t_{nj}) + \mathcal{E}_{nj}, \quad 1 \le j \le n, \quad n \in \mathbb{N},$$
(5)

where  $\mathcal{E}_{n1}, \dots, \mathcal{E}_{nn}$  are independent and identically distributed (iid) random variables with  $E(\mathcal{E}_{nj}) = 0$ ,  $Var(\mathcal{E}_{ni}) = \sigma^2 = 1$ .

Since our results keep true when  $\sigma^2$  is replaced by a consistent estimator, we can put  $\sigma^2 = 1$  without loss of generality.

Putting  $\tau_n := (t_{n1}, \dots, t_{nn})$  we define

$$Y(\tau_n) \coloneqq (Y(t_{n1}), \dots, Y(t_{nn}))^T,$$
  

$$g(\tau_n) \coloneqq (g(t_{n1}), \dots, g(t_{nn}))^T,$$
  

$$f_i(\tau_n) \coloneqq (f_i(t_{n1}), \dots, f_i(t_{nn}))^T, \quad i = 1, \dots, d$$
  

$$\varepsilon_n \coloneqq (\varepsilon_{n1}, \dots, \varepsilon_{nn})^T,$$

then the triangular array of observations can be written by  $Y(\tau_n) = g(\tau_n) + \varepsilon_n$ ,  $n \in \mathbb{N}$ .

Let us consider the vector spaces  $W := [f_1(\cdot), \dots, f_d(\cdot)]$  and  $W_{\tau_n} := [f_1(\tau_n), \dots, f_d(\tau_n)]$ spanned by the functions  $f_1(\cdot), \dots, f_d(\cdot)$  and by the vectors  $f_1(\tau_n), \dots, f_d(\tau_n)$ , respectively. The intrinsic test problem is

$$H_0: g(\cdot) \in W \coloneqq [f_1(\cdot), \dots, f_d(\cdot)] \quad \text{against} \quad K: g(\cdot) \notin W.$$
(6)

Without assumptions on the smoothness of the unknown true regression function g we cannot decide the above test problem by a finite sample. Hence, without any assumption the above test problem is not well stated. We assume:

$$g \in BV[0,1]$$

i.e. g has bounded variation.

Note, that this assumption is no restriction in practice. Hence, the known regression functions  $f_1, \ldots, f_d$  must belong also to BV[0,1]. Moreover, we assume

$$f_1, \dots, f_d \in BV[0,1] \cap C[0,1], \tag{8}$$

where C[0,1] is the set of continuous functions on [0,1]. By our assumption on g we know that there is an  $n^* \in \mathbb{N}$  with  $g(\tau_n) \notin W_{\tau_n}$  for all  $n > n^*$  if  $g(\cdot) \notin W$ . Therefore, if n is large enough, the original test problem can be decided by a test for  $H_0: g(\tau_n) \in W_{\tau_n} := [f_1(\tau_n), \dots, f_d(\tau_n)]$  against  $K: g(\tau_n) \notin W_{\tau_n}$ .

Under the null hypothesis we have the linear model

$$Y(\tau_n) = f^T(\tau_n)\beta + \varepsilon_n, \quad n \in \mathbb{N},$$
(9)

where  $f^T(\tau_n) = (f_1(\tau_n), \dots, f_d(\tau_n)) =: X_n$  is the design matrix. To state our results it is convenient to define the partial sums operator:  $T_n : \mathbb{R}^n \to C[0,1], \ \mathbf{a} = (a_1, \dots, a_n)^T \mapsto T_n(\mathbf{a})(z), \ z \in [0,1],$ where

$$T_n(\mathbf{a})(z) = \sum_{i=1}^{|nz|} a_i + (nz - [nz])a_{[nz]+1}, z \in [0,1].$$

Here we used  $[s] = \max\{n \in \mathbb{N}_0 | n \le s\}$  and  $\sum_{i=1}^0 a_i = 0$ . Let us define  $b_i = a_1 + \ldots + a_i$ ,  $i = 1, \ldots, n$ , then  $T_n(a), a = (a_1, \ldots, a_n)^T$ , is shown in Figure 2.



Figure 2. The function  $T_n(a)(\cdot)$ 

The stochastic process  $\frac{1}{\sqrt{n}}T_n(\mathcal{E}_n)$  converges weakly to the Brownian motion *B* for  $n \to \infty$  by Donsker's Theorem.

#### 3. Residual Partial Sum Processes

Let  $\operatorname{pr}_{W_{\tau_n}}$  be the orthogonal projector onto  $W_{\tau_n}$  in  $\mathbb{R}^n$  with respect to the Euclidean inner product. Thus the least squares estimator for  $X_n\beta$  is given by  $pr_{W_{\tau_n}}Y(\tau_n)$  and the vector of the corresponding least squares residuals by  $r_n \coloneqq Y(\tau_n) - pr_{W_{\tau_n}}Y(\tau_n)$ .

If the null-hypothesis  $H_0$  is true, then we have

$$r_n = Y(\tau_n) - pr_{W_{\tau_n}}Y(\tau_n) = \sum_{i=1}^d f_i(\tau_n)\beta_i + \varepsilon_n - pr_{W_{\tau_n}}(\sum_{i=1}^d f_i(\tau_n)\beta_i + \varepsilon_n) = \varepsilon_n - pr_{W_{\tau_n}}\varepsilon_n$$

We are interested in:  $\frac{1}{\sqrt{n}}T_n(r_n) = \frac{1}{\sqrt{n}}T_n(\varepsilon_n - pr_{W_{\tau_n}}\varepsilon_n)$  for  $n \to \infty$  to be able to state a test with the help of the limit distribution. Note that we have not assumed any distribution for the error. To be able to state the following theoretical main result we need some further notation. Let  $L^2(\lambda)$  be the space of square integrable functions with respect to the Lebesgue measure  $\lambda$ . Then for each  $h \in L^2(\lambda)$  the function  $s_h(\cdot) \coloneqq \int_{[0,\cdot]} h d\lambda$  is called signal. Let  $H = \{s_k \mid h \in L_2(\lambda)\}$  be furnished with the inner product  $\langle s_{h_1}, s_{h_2} \rangle_H = \int_{[0,1]} h_1 h_2 d\lambda = \langle h_1, h_2 \rangle_{L^2(\lambda)}$ . Then we consider the linear subspace  $W_{\rm H} \coloneqq [s_{f_1}(\cdot), \dots, s_{f_4}(\cdot)]$  in H

and the projection  $pr_{W_H}$  onto  $W_H$  in H which can be extended to a projection defined on C[0,1].

#### Theorem 1. (MacNeill [20, 21], Bischoff [2]).

Let  $f_1, \ldots, f_d \in BV[0,1] \cap C[0,1]$ . Then under the null hypothesis ' $H_0 : g \in W = [f_1, \ldots, f_d]$ ' we have  $\frac{1}{\sqrt{n}}T_n(Y(\tau_n) - pr_{W_{\tau_n}}Y(\tau_n))$  converges weakly to  $B - pr_{W_H}B$ . Moreover, let  $g \in BV[0,1] \setminus W$ , then for the sequence of alternatives  $\frac{1}{\sqrt{n}}g(\tau_n), n \in \mathbb{N}$ , we have

$$\frac{1}{\sqrt{n}}T_{n}\left(\frac{1}{\sqrt{n}}g\left(\tau_{n}\right)+\varepsilon_{n}-pr_{W_{\tau_{n}}}\left(\frac{1}{\sqrt{n}}g\left(\tau_{n}\right)+\varepsilon_{n}\right)\right) \text{ converges weakly to}$$

$$s_{g}-pr_{W_{H}}s_{g}+B-pr_{W_{H}}B, \text{ where } s_{g}-pr_{W_{H}}s_{g}\neq 0.$$

The test problem stated in (6) can be transformed to the equivalent test problem

$$H_0: s_g(\cdot) \in W_H := [s_{f_1}(\cdot), \dots, s_{f_d}(\cdot)] \text{ against } K: s_g(\cdot) \notin W_H.$$

Next, we consider the problem and the data discussed in Section 1. There we have

$$f(t) = f_{1}(t) = \mathbf{1}_{[0,1]}(t), t \in [0,1], W = [\mathbf{1}_{[0,1]}(\cdot)], W_{\tau_{n}} = [\mathbf{1}_{[0,1]}(\tau_{n})],$$
  

$$s_{f}(z) = s_{f_{1}}(z) = \int_{[0,z]} \mathbf{1}_{[0,1]}(t) \quad \lambda(dt) = z, \quad z \in [0,1], \quad W_{H} = [s_{f}(\cdot)] = [\operatorname{id}_{[0,1]}(\cdot)].$$

In this case we get the following limit process for the residual partial sum process under the null hypothesis  $g \in W$ :

$$B(z) - \operatorname{pr}_{W_H} B(z) = B(z) - B(1)z =: B_0(z), z \in [0,1],$$

where  $B_0$  is the Brownian bridge.

Under the alternative  $K: g \notin W$ , however, we have to scale g with  $\frac{1}{\sqrt{n}}$  to get  $s_g - pr_{W_H}s_g + B_0$  as limit of the residual partial sum process with a signal  $s_g - pr_{W_H}s_g \neq 0$ . Let  $r_n$  be the vector of the least squares residuals of the data discussed in Section 1 and shown on Figure 1. Then the path of the residual partial sums process  $\frac{1}{\hat{\sigma}\sqrt{n}}T_n(r_n)$  is given in Figure 3. Here  $\hat{\sigma}$  is a consistent estimator for  $\sigma$ .

#### 4. Tests of Kolmogorov–Smirnov Type

Let us consider our data again to check whether the quality is constant. Hence, we want to test by a size  $\alpha$  test,  $\alpha \in (0,1)$ ,

$$H_0: g \in [\mathbf{1}_{[0,1]}(\cdot)] = W, \quad \text{against} \quad K: g \notin W.$$

$$\tag{10}$$

For that we take a two-sided test of Kolmogorov-Smirnov type based on the residual partial sum process:

Reject 
$$H_0 \Leftrightarrow \exists z \in [0,1]: \frac{1}{\sqrt{n}} T_n(Y(\tau_n) - pr_{W_n}Y(\tau_n))(z) \notin [\ell(z), u(z)],$$
 (11)

where  $\ell, u: [0,1] \to \mathbb{R}$  are chosen in such a way that

$$P(\exists z \in [0,1]: B_0(z) \notin [\ell(z), u(z)]) = \alpha.$$

$$(12)$$

Thus (11) is an asymptotic size  $\alpha$  test for (10) by Theorem 1.

In many cases one is even more interested whether the quality is getting worse. For the quality problem we know that the greater the value the better the quality. Hence, we want to test by a size  $\alpha$  test,  $\alpha \in (0,1)$ ,  $H_0: g$  is constant against  $K: g \in BV[0,1]$  is decreasing and is not constant.

For that we consider a one-sided test of Kolmogorov type based on the residual partial sum process:

Reject 
$$H_0 \Leftrightarrow \exists z \in [0,1]: \frac{1}{\sqrt{n}} T_n(Y(\tau_n) - pr_{W_n}Y(\tau_n))(z) \ge u(z),$$
 (13)

where  $u \in [0,1] \to \mathbb{R}$  is chosen in such a way that  $P(\exists z \in [0,1]: B_0(z) \ge u(z)) = \alpha$ .

Thus (13) is an asymptotic size  $\alpha$  test for (12) by Theorem 1. Note that the above form of the test is adequate for the problem to test for decreasing quality because in that case the mean of the residuals are decreasing. Then the signal  $s - pr_{W_H}s$  of the residual partial sum process is a concave function with  $(s - pr_{W_H}s)(0) = (s - pr_{W_H}s)(1) = 0$ , where  $W_H = [\operatorname{id}_{[0,1]}]$ . For the data of our quality problem we used a Kolmogorov test with constant boundary *u*. On Figure 3 is shown the residual partial sum process of the data together with the constant boundary *u* which corresponds to the size  $\alpha = 0.01$ .



Figure 3.  $\frac{1}{\hat{\sigma}\sqrt{n}}T_n(r_n)$  for the data given in Figure 1 with the constant boundary *u* of the Kolmogorov test

for the size  $\alpha = 0.01$ 

Since  $\frac{1}{\hat{\sigma}\sqrt{n}}T_n(r_n)$  crosses the boundary *u* for the size  $\alpha = 0.01$  we reject the null hypothesis at this significance level  $\alpha$ .

#### 5. Further Results and Generalisations

Finally, we cite some papers that considered related problems. In Bischoff [1] residual partial sum processes for not necessarily equidistant design points are established. Tang and MacNeill [25] determined residual partial sum processes for time series. Xie and MacNeill [27], Bischoff and Somayasa [11] investigated residual partial sum processes for a linear regression model with a multivariate experimental region. Bischoff and Gegg [3] considered residual partial sum processes with multivariate response. The power of the Kolmogorov(–Smirnov) test and the corresponding problem of boundary crossing probabilities of Gaussian processes is investigated in Bischoff et al. [6–9], Bischoff and Hashorva [4].

#### References

- Bischoff, W. A functional central limit theorem for regression models, *Ann. Statist.*, Vol. 26, 1998, pp. 1398–1410.
- 2. Bischoff, W. The structure of residual partial sums limit processes of linear regression models, *Theory of Stochastic Processes*, Vol. 8 (24), No 1–2, 2002, pp. 23–28.
- 3. Bischoff, W. and Gegg, A. Partial sums process to check regression models with multiple correlated response with an application: test for a change-point in profile data, *J. Multivariate Anal.*, 2010. (Submitted)
- 4. Bischoff, W., Hashorva, E. A lower bound for boundary crossing probabilities of Brownian bridge with trend, *Statist. Probab. Lett.*, Vol. 74, 2005, pp. 265–271.
- 5. Bischoff, W., Hashorva, E. and Hüsler, J. An asymptotic result for boundary Crossing probabilities of Brownian motion with trend, *Communications in Statistics: Theory and Method*, Vol. 36:16, 2007, pp. 2821–2828.
- 6. Bischoff, W., Hashorva, E., Hüsler, J. and Miller, F. Asymptotics of a boundary crossing probability of a Brownian bridge with general trend, *Methodology and Computing in Applied Probability*, Vol. 5, 2003a, pp. 271–287.
- Bischoff, W., Hashorva, E., Hüsler, J. and Miller, F. Exact asymptotics for boundary crossings of the Brownian bridge with trend with application to the Kolmogorov test, *Ann. Inst. Statist. Math.*, Vol. 55, 2003b, pp. 849–864.

- 8. Bischoff, W., Hashorva, E., Hüsler, J. and Miller, F. On the power of the Kolmogorov test to detect the trend of a Brownian bridge with applications to a change-point problem in regression models. *Stat. Probab. Lett.*, Vol. 66, 2004, pp. 105–115.
- 9. Bischoff, W., Hashorva, E., Hüsler, J. and Miller, F. Analysis of a change-point regression problem in quality control by partial sums processes and Kolmogorov type tests, *Metrika*, Vol. 62, No 1, 2005, pp. 85–98.
- Bischoff, W. and Miller, F. Asymptotically optimal tests and optimal designs for testing the mean in regression models with applications to change-point problems, *Ann. Inst. Statist. Math.*, Vol. 52, 2000, pp. 658–679.
- 11. Bischoff, W. and Somayasa, W. The limit of the partial sums process of spatial least squares residuals, *J. Multivariate Anal.*, Vol. 100, No 10, 2009, pp. 2167–2177.
- 12. Brown, R. L., Durbin, J. and Evans, J. M. Techniques for testing the constancy of regression relationships over time, *J. Roy. Stat. Soc.*, Vol. B 37, 1975, pp. 149–192.
- 13. Gardner, L. A. On detecting changes in the mean of normal variates, *Ann. Math. Stat.*, Vol. 40, 1969, pp. 16–126.
- 14. Jandhyala, V. K. and Al-Saleh, J. A. Parameter changes at unknown times in non-linear regression, *Environmetrics*, Vol. 10, 1999, pp. 711–724.
- Jandhyala, V. K. and MacNeill, I. B. The residual process for non-linear regression, *J. Appl. Prob.*, Vol. 22, 1985, pp. 957–963.
- Jandhyala, V. K. and MacNeill, I. B. Residual partial sum limit processes for regression models with applications to detecting parameter changes at unknown times, *Stoch. Process. Appl.*, Vol. 33, 1989, pp. 309–323.
- 17. Jandhyala, V. K. and MacNeill, I. B. Tests for parameter changes at unknown times in linear regression models, *J. Statist. Plan. Inf.*, Vol. 27, 1991, pp. 291–316.
- 18. Jandhyala, V. K. and MacNeill, I. B. Iterated partial sum sequences of regression residuals and tests for change-points with continuity constraints, *J.R. Statist. Soc.*, Vol. 59, 1999, pp. 147–156.
- 19. Jandhyala, V. K., Zacks, S. and El-Shaarawi, A. H. Change-point methods and their applications: Contributions of Ian MacNeill, *Environmetrics*, Vol. 10, 1999, pp. 657–676.
- 20. MacNeill, I. B. Property of sequences of partial sums of polynomial regression residuals with applications to tests for change of regression at unknown times, *Ann. Statist.*, Vol. 6, 1978a, pp. 422–433.
- 21. MacNeill, I. B. Limit processes for sequences of partial sums of regression residuals, *Ann. Prob.*, Vol. 6, 1978b, pp. 695–698.
- 22. Sen, P. K. Invariance principles for recursive residuals, Ann. Statist., Vol. 10, 1982, pp. 307-312.
- 23. Sen, A. and Srivastava, M. S. On tests for detecting change in mean when variance is unknown, *Ann. Inst. Stat. Math.*, Vol. 27, 1975, pp. 479–486.
- 24. Stute, W. Nonparametric model checks for regression, Ann. Statist., Vol. 25, No 2, 1997, pp. 613-641.
- 25. Tang, S. M. and MacNeill, I. B. The effect of serial correlation on tests for parameter change at unknown time, *Ann. Stat.*, Vol. 21, No 1, 1993, pp. 552–575.
- 26. Watson, G. S. Detecting a change in the intercept in multiple regression, *Stat. Probab. Lett.*, Vol. 23, 1995, pp. 69–72.
- 27. Xie, L. and MacNeill, I. B. Spatial residual processes and boundary detection, *South African Statist. J.*, Vol. 40, No 1, 2006, pp. 33–53.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 14–18 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# BETA-DISTRIBUTION MODELS IN STOCHASTIC PROJECT MANAGEMENT

**D.** Greenberg<sup>1</sup>, A. Ben-Yair<sup>2</sup>

<sup>1</sup>Department of Economics and Business Administration, Faculty of Social Science Ariel University Center (AUC) of Samaria P.O. Box 3, Ariel, 40700, Israel E-mail: dorongreen2@gmail.com <sup>2</sup>The Department of Industrial Engineering and Management, SCE-Shamoon College of Engineering

> Beer-Sheva 84100, Israel E-mail: avnerb@sce.ac.il

A research is undertaken to justify the use of beta-distribution p.d.f. for man-machine type activities under random disturbances. The case of using one processor, i.e., a single resource unit, is examined. It can be proven theoretically that under certain realistic assumptions the random activity – time distribution satisfies the beta p.d.f. Changing more or less the implemented assumptions, we may alter to a certain extent the structure of the p.d.f. At the same time, its essential features (e.g. asymmetry, unimodality, etc.) remain unchanged. The outlined above research can be applied to semi-automated activities, where the presence of man-machine influence under random disturbances is, indeed, very essential. Those activities are likely to be considered in organization systems (e.g. in project management), but not in fully automated plants.

**Keywords:** random activity duration, time – activity beta-distribution, operating by means of a single processor, convergence to a beta-distribution "family"

#### 1. Introduction

In PERT analysis [1–24, etc.] the activity-time distribution is assumed to be a beta-distribution, and the mean value and variance of the activity time are estimated on the basis of the "optimistic", "most likely" and "pessimistic" completion times, which are subjectively determined by an analyst. The creators of PERT [3, 17] worked out the basic concepts of PERT analysis, and suggested the estimates of the mean and variance values

$$\mu = \frac{1}{6} \left( \mathbf{a} + 4\mathbf{m} + \mathbf{b} \right),\tag{1}$$

$$\sigma^2 = \frac{1}{36} (b-a)^2,$$
 (2)

subject to the assumption that the probability density function (p.d.f.) of the activity time is

$$f_{y}(y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{(y-a)^{\alpha-1}(b-y)^{\beta-1}}{(b-a)^{\alpha+\beta-1}}, \quad a < y < b, \quad \alpha, \beta > 0.$$
(3)

Here a is the optimistic time, b – the pessimistic time, and m stands for the most likely (modal) time.

Since in PERT applications parameters a and b of p.d.f. (3) are either known or subjectively determined, we can always transform the density function to a standard form,

$$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha - 1} (1 - x)^{\beta - 1}, \quad 0 < x < 1, \quad \alpha, \beta > 0,$$
(4)

where  $x = \frac{y-a}{b-a}$  has the following parameters:

$$\mu_{x} = \frac{\mu_{y} - a}{b - a}, \quad \sigma_{x} = \frac{\sigma_{y}}{b - a}, \quad m_{x} = \frac{m_{y} - a}{b - a}.$$
(5)

Let  $\alpha - 1 = p$ ,  $\beta - 1 = q$ . Then p.d.f. (4) becomes

$$f(x) = \frac{\Gamma(p+q+2)}{\Gamma(p+1)\Gamma(q+1)} x^{p}(1-x)^{q}, \quad 0 < x < 1, \quad p,q > -1,$$
(6)

with the mean, variance and mode as follows:

$$\mu_{\rm x} = \frac{{\rm p} + {\rm I}}{{\rm p} + {\rm q} + 2},\tag{7}$$

$$\sigma_{\rm x}^2 = \frac{({\rm p}+1)({\rm q}+1)}{{\rm p}+{\rm q}+2},\tag{8}$$

$$m_x = \frac{p}{p+q}.$$
(9)

From (6) and (9) it can be obtained

$$f(x) = \frac{\Gamma(p+q+2)}{\Gamma(p+1)\Gamma(q+1)} x^{p}(1-x)^{p(l/m_{x}-1)}.$$
(10)

Thus, value  $m_x$ , being obtained from the analyst's subjective knowledge, indicates the density function. On the basis of statistical analysis and some other intuitive arguments, the creators of PERT assumed that  $p+q \cong 4$ . It is from that assertion that estimates (1) and (2) were finally obtained, according to (6–9).

Although the basic concepts of PERT analysis have been worked out many years ago [3, 17], they are open till now to considerable criticism. Numerous attempts have been made to improve the main PERT assumptions for calculating the mean  $\mu_x$  and variance  $\sigma_x^2$  of the activity-time on the basis of the analyst's subjective estimates. In recent years, a very sharp discussion [7, 10, 14, 21] has taken place in order to raise the level of theoretical justifications for estimates (1) and (2).

Grubbs [12] pointed out the lack of theoretical justification and the unavoidable defects of the PERT statements, since estimates (1) and (2) are, indeed, "rough" and cannot be obtained from (3) on the basis of values a, m and b determined by the analyst. Moder [18–19] noted that there is a tendency to choose the most likely activity – time m much closer to the optimistic value a than to the pessimistic one, b, since the latter is usually difficult to determine and thus is taken conservatively large. Moreover, it is shown [8] that value m, being subjectively determined, has approximately one and the same relative location point in [a,b] for different activities. This provides an opportunity to simplify the PERT analysis at the expense of some additional assumptions. McCrimmon and Ryavec [16], Lukaszewicz [15] and Welsh [22] examined various errors introduced by the PERT assumptions, and came to the conclusion that these errors may be as great as 33%. Murray [20] and Donaldson [4] suggested some modifications of the PERT analysis, but the main contradictions nevertheless remained. Farnum and Stanton [6] presented an interesting improvement of estimates (1) and (2) for cases when the modal value m is close to the upper or lower limits of the distribution. This modification, however, makes the distribution law rather uncertain, and causes substantial difficulties to simulate the activity network.

In this paper, a research will be undertaken to develop some theoretical justifications for using the beta-distribution p.d.f.

#### 2. The Operation's Description

We will consider a man-machine operation which is carried out by one processor, i.e., by one resource unit. The processor may be a machine, a proving ground, a department in a design office, etc.

Assume that the operation starts to be processed at a pregiven moment  $T_0$ . The completion moment F of the operation is a random value with distribution range  $[T_1, T_2]$ . Moment  $T_1$  is the operation's

completion moment on condition that the operation will be processed without breaks and without delays, i.e., value  $T_I$  is a pregiven deterministic value. Assume, further, that the interval  $[T_0, T_1]$  is subdivided into n equal elementary periods with length  $(T_1 - T_0)/n$ . If within the first elementary period  $[T_0, T_0 + (T_1 - T_0)/n]$  a break occurs, it causes a delay of length  $\Delta = (T_2 - T_1)/n$ . The operation stops to be processed within the period of delay in order to undertake necessary refinements, and later on proceeds functioning with the finishing time of the first elementary period  $T_0 + (T_1 - T_0)/n + (T_2 - T_1)/n = T_0 + (T_2 - T_0)/n$ .

It is assumed that there cannot be more than one break in each elementary period. The probability of a break at the very beginning of the operation is set to be p. However, in the course of carrying out the operation, the latter possesses certain features of self-adaptivity, as follows:

- the occurrence of a break within a certain elementary period results in increasing the probability of a new break at the next period by value η, and
- on the contrary, the absence of a break within a certain period decreases the probability of a new break within the next period, practically by the same value.

#### 3. The Concept of Self-Adaptivity

The probabilistic self-adaptivity can be formalized as follows:

Denote  $A_i^k$  the event of occurrence of a break within the (i+1)-th elementary period, on condition, that within the i preceding elementary periods k breaks occurred,  $1 \le k \le i \le n$ . It is assumed that relation

$$P\left(A_{i}^{k}\right) = \frac{p + k \cdot \eta}{1 + i \cdot \eta}$$

$$\tag{11}$$

holds. Note that (11) is, indeed, a realistic assumption.

Relation (11) enables obtaining an important assertion. Let  $P(A_i^0)$  be the probability of the occurrence of a break within the (i+1)-th period on condition, that there have been no breaks at all as yet. Since

$$P(A_i^0) = \frac{p}{1+i \cdot p},$$
(12)

it can be well-recognized that relation

$$\frac{P(A_i^{k+1}) - P(A_i^k)}{P(A_i^0)} = \frac{\eta}{p}$$
(13)

holds. Thus, an assertion can be formulated as follows:

<u>Assertion</u>. Self-adaptivity (11) results in a probability law for delays with a constant ratio (13) for a single delay.

#### 4. Calculating the Activity-Time Distribution

Let us calculate the probability  $P_{m,n}$  of obtaining *m* delays within *n* elementary periods, i.e., the probability of completing the operation at the moment  $F = T_1 + m \cdot \Delta = T_1 + \frac{m}{n} (T_2 - T_1)$ .

The number of sequences of n elements with m delays within the period  $[T_0, F]$  is equal  $C_n^m$ , while the probability of each such sequence equals

$$\frac{\left[\prod_{i=0}^{m-1} (p+i\eta)\right] \left[\prod_{i=0}^{n-m-1} (1-\eta+i\eta)\right]}{\prod_{i=0}^{n-1} (1+i\eta)} .$$
(14)

Relation (14) stems from the fact that if breaks occurred within h periods and did not occur within k periods, the probability of the occurrence of the delay at the next period is equal

$$\frac{\mathbf{p} + \mathbf{h}\eta}{1 + (\mathbf{k} + \mathbf{h})\eta} , \qquad (15)$$

while the probability of the delay's non-appearance at the next period satisfies

$$\frac{1-\eta+k\eta}{1+(k+h)\eta} .$$
(16)

Using (14-16), we finally obtain

$$P_{m,n} = C_n^m \frac{\left[\prod_{i=0}^{m-1} (p+i\eta)\right] \left[\prod_{i=0}^{n-m-1} (1-\eta+k\eta)\right]}{\prod_{i=0}^{n-1} (1+i\eta)} .$$
(17)

Note that  $\eta = 0$ , i.e., the absence of self-adaptivity, results in a regular binomial distribution. Let us now obtain the limit value  $P_{m,n}$  on condition that  $n \to \infty$ . From relation (17) we obtain

$$\frac{P_{m+1,n}}{P_{m,n}} = \frac{n-m}{m+1} \frac{p+m\eta}{1-p+(n-m-1)\eta} .$$

$$Denoting \quad \frac{p}{\eta} = \alpha , \quad \frac{p}{\eta} \left(\frac{1}{p}-1\right) = \beta , \text{ we obtain}$$

$$\frac{P_{m+1,n}-P_{m,n}}{P_{m,n}} = \frac{(\alpha-1)n+(2-\alpha-\beta)m-\beta+1}{(m+1)(\beta+n-m-1)} = \frac{(\alpha-1)+(2-\alpha-\beta)\frac{m}{n}+\frac{1-\beta}{n}}{n\frac{m+1}{n}\left(1-\frac{m+1}{n}+\frac{\beta}{n}\right)} .$$
(18)

Denoting m/n = x,  $(m+1)/n = x + \Delta x$ ,  $P_{m,n} = y$ ,  $P_{m+1,n} = y + \Delta y$ , via convergence  $n \to \infty$ or  $\Delta x \to 0$  and, later on, by means of integration, we finally obtain  $y = C x^{\alpha-1} (1-x)^{\beta-1}$ . (19)

It can be well-recognized that the p.d.f. of random value  $\xi = \lim_{n \to \infty} \frac{m}{n}$  satisfies

$$p_{\xi}(x) = \frac{1}{B(\alpha,\beta)} x^{\alpha-1} (1-x)^{\beta-1} , \qquad (20)$$

where  $B(\alpha, \beta)$  represents the Eiler's function. Thus, relation (20) practically coincides with (10).

Thus,  $\xi$  is a random value with the beta-distribution activity – time p.d.f. By transforming x = (y-a)/(b-a), we obtain the well-known p.d.f. (3).

#### Conclusions

The following conclusions can be drawn from the study:

• Under certain realistic assumptions we have proven theoretically that the activity-time distribution satisfies the beta-distribution with p.d.f. (3) being used in PERT analysis.

- Changing more or less the implemented assumptions, we may alter to a certain extent the structure of the p.d.f. At the same time, *its essential features* (e.g. asymmetry, unimodality, etc.) remain unchanged.
- The outlined above research can be applied to semi-automated activities, where the presence of manmachine influence under random disturbances is, indeed, very essential. Those activities are likely to be considered in organization systems (e.g. in project management), but not in fully automated plants.

#### References

- 1. Battersby, A. Network Analysis for Planning and Scheduling, 3<sup>rd</sup> edition. London: MacMillan, 1970.
- 2. Berny, J. A new distribution functions for risk analysis, *J. Oper. Res. Soc.*, Vol. 40, No 12, 1989, pp. 1121–1127.
- 3. Clark, C. E. The PERT model for the distribution of an activity, Opns. Res., Vol. 10, 1962, pp. 405–406.
- 4. Donaldson, W. A. The estimation of the mean and variance of PERT activity time, *Opns. Res.*, Vol. 13, 1965, pp. 382–385.
- 5. Elmaghraby, S. E. Activity Networks: Project Planning and Control by Network Models. New-York: Wiley, 1977.
- 6. Farnum, N. R. and L. W. Stanton. Some results concerning the estimation of beta distribution parameters in PERT, *J. Oper. Res. Soc.*, Vol. 38, 1987, pp. 287–290.
- 7. Gallagher, C. A note on PERT assumptions, Mgmt. Sci., Vol. 33, 1987, pp. 1360–1362.
- 8. Golenko-Ginzburg, D. *Statistical Models in Network Planning and Control*, Moscow: Nauka, 1966. (In Russian)
- 9. Golenko-Ginzburg, D. On the distribution of activity time in PERT, *J. Oper. Res. Soc.*, Vol. 39, No 8, 1988, pp. 767–771.
- 10. Golenko-Ginzburg, D. PERT assumptions revisited, Omega, Vol. 17, No 4, 1989, pp. 393-396.
- 11. Gonik, A. *Planning and Controlling Multilevel Stochastic Projects: Ph.D. Thesis.* Beer-Sheva: Ben-Gurion University of the Negev, 1995.
- Grubbs, F. E. Attempts to validate certain PERT statistics or "picking on PERT", Opns. Res., Vol. 10, 1962, pp. 912–915.
- 13. Kelley, J. E., Jr. Critical path planning and scheduling: Mathematical basis, *Opns. Res.*, Vol. 9, No 3, 1961, pp. 296–320.
- 14. Littlefield, T. K., Jr. and P. H. Randolph. Another note on PERT times, *Mgmt. Sci.*, Vol. 33, 1987, pp. 1357–1359.
- 15. Lukaszevicz, J. On the estimation of errors introduced by standard assumptions concerning the distribution of activity duration PERT calculations, *Opns. Res.*, Vol. 13, 1965, pp. 326–327.
- MacCrimmon, K. R. and C. A. Ryaveck. An analytical study of the PERT assumptions, *Opns. Res.*, Vol. 12, 1964, pp. 16–37.
- 17. Malcolm, D., Roseboom, J., Clark, C. and W. Fazar. Application of a technique for research and development program evaluation, *Opns. Res.*, Vol. 7, 1959, pp. 646–669.
- 18. Moder, J. J., Phillips, C. R., Davis, E. W. *Project Management with CPM and PERT and Precedence Diagramming*. New-York: Van-Nostrand Reinhold Co., Inc., 1983.
- 19. Moder, J. J., Cecil, R. P. Project Management with CPM and PERT. New-York: Van-Nostrand Reinhold Co., Inc., 1970.
- 20. Murray, J. E. *Consideration of PERT Assumptions*. Ann Arbour, Michigan: Conduction Corporation, 1962.
- 21. Sasieni, M. W. A note on PERT times, Mgmt. Sci., Vol. 32, 1986, pp. 1652–1653.
- 22. Welsh, D. Errors introduced by a PERT assumption, Opns. Res., Vol. 13, 1965, pp. 141-143.
- 23. Williams, T. M. Practical use of distributions in network analysis, J. Oper. Res. Soc., Vol. 43, No 3, 1992, pp. 265–270.
- 24. Williams, T. M. What are PERT estimates? J. Oper. Res. Soc., Vol. 46, No 12, 1995, pp. 1498–1504.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 19–30 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# CONTINUOUS MODELS OF CURRENT STOCK OF DIVISIBLE PRODUCTIONS

E. Kopytov<sup>1</sup>, S. Guseynov<sup>1, 2</sup>, E. Puzinkevich<sup>1</sup>, L. Greenglaz<sup>1</sup>

<sup>1</sup>Transport and Telecommunication Institute Lomonosova Str. 1, Riga LV-1019, Latvia E-mail: kopitov@tsi.lv / edvins.puzinkevics@du.lv

<sup>2</sup>Institute of Mathematical Sciences and Information Technologies, University of Liepaja Liela Str. 14, Liepaja LV-3401, Latvia E-mail: sh.e.guseinov@inbox.lv

In the given paper we investigate the problem of constructing continuous and unsteady mathematical models to determine the volumes of current stock of divisible productions using apparatus and equations of mathematical physics. It is assumed that time of production distribution and replenishment is continuous. The constructed models are stochastic, and have different levels of complexity, adequacy and application potentials. The simple models are constructed using the theory of ordinary differential equations, for construction of more complex models the theory of partial differential equations is applied. Furthermore for some of proposed models we have found an analytical solution in the closed form, and for some of proposed models the discretization is carried out using stable difference schemes.

Keywords: inventory control model, current stock, divisible production, equations of mathematical physics

#### 1. Introduction

One of the central problems of the inventory control theory is to find an optimal or quasi optimal solution to the task of ordering products to be supplied. Of no less interest it is the task of determining the current stock of certain products (sold by the piece or indivisible products and dry or divisible products) at any given moment of a fixed time span, with any random factors taken into account. By "current stock" we denote the quantity (volume) of the product accumulated in the stock, which is used for regular distribution/replenishment. Quite a lot of different types of models of varying complexity, purpose and adequacy have been developed in the inventory control theory. Most of the existing mathematical models in this theory consider indivisible products (for example, see [1-3]). We can classify these models taking in account their different properties: deterministic and stochastic, linear and nonlinear, single- and multi-product, discrete and continuous models, etc. [3].

The present paper studies construction of continuous and unsteady mathematical models for calculating the volume of current stock of divisible production "from scratch" using apparatus and equations of mathematical physics [4]. The suggested models are stochastic ones and have different levels of complexity, adequacy and application potentials. The simple models are constructed using the theory of ordinary differential equations, for construction of more complex models the theory of partial differential equations is applied.

# 2. Construction of the Continuous Stochastic Model for Determining the Volume of the Current Stock of Homogeneous Divisible Production

In the present section we construct the continuous stochastic mathematical model for determining the volume of current stock of divisible production. For this purpose, we will use the apparatus of mathematical physics and the continuum principle (for example, see [4]); as modelling language will be chosen language of ordinary differential equations (ODE). Before introducing the simplifying assumptions, which are required for modelling, as well as variables, parameters and functions that are describing and coupling the initial data of the simulated process with unknown quantities of the current stock dynamics, we will consider briefly the issue of stochasticity of the mathematical model under construction. Namely, to construct the stochastic model, we can proceed in the following two ways:

- the current stock to be determined is not supposed to be an accidental quantity, but after the introduction of a change rate the constructed model is supplied with all random factors which visibly influence unknown rate of the current stock change. In this case the obtained relation (in the form of the above

mentioned ODE) with regard to unknown volume of the current stock and rate of its change is a functional relationship among unknown volume, rate of its change, and accidental factors influencing the current stock dynamics. In other words, in the obtained model, unknown volume, which initially did not seem to be assumed as an accidental value (stochastic value of a random function, to be more specific), due to the obtained ODE and corresponding conditions appears dependent on the random quantities taken into account, i.e. unknown volume of the current stock is a function of the accidental quantities;

- the current stock is initially taken to be a random quantity, and this suggestion is taken into account when constructing the model.

The first of these ways is selected for the description of the mathematical model that will follow. It is worth mentioning in the way of a preliminary note that this choice will result in the construction of a stochastic model represented by the Ito-type differential equation (for example, see [5] as well as works [6], [7], [8] on stochastic differential equations).

Now we can start constructing the mathematical model "from the scratch". Let us assume that the current stock volume of the considered homogeneous divisible production at the moment t equals to x(t). It is required that  $x(t) \in C[T_s, T_e]$ ;  $\exists x'(t) \forall t \in [T_s, T_e]$ , where  $[T_s, T_e]$  is a segment of time during which the dynamics of the current stock change is being studied, by  $T_s$  and  $T_e$  we denote the initial and final moments of this period of time, respectively. The requirement  $x(t) \in C[T_s, T_e]$  is easy to interpret economically, and it is met if we assume that the current stock x(t) is being constantly distributed/replenished. The requirement  $\exists x'(t) \forall t \in [T_s, T_e]$  is a purely mathematical one, i.e. it is necessary to ensure a mathematical correctness of the model.

If an increase of the current stock volume x(t) is as

$$\Delta x(t) \stackrel{\text{def}}{=} x(t + \Delta t) - x(t), \ \Delta t > 0, \ t + \Delta t \le T_e, \text{ then}$$

$$\frac{dx(t)}{dt} \equiv \lim_{\Delta t \to 0} \frac{\Delta x(t)}{\Delta t}, \tag{1}$$

and this quantity designates the change rate of the current stock volume at a given time t.

The rate  $\frac{dx(t)}{dt}$  derived from (1) is completely analogous to the rate of a material point of continuous medium moving in metric space. It is then useful to find out the factors or reasons causing the change x(t) and, consequently, trigger the existence of  $\frac{dx(t)}{dt}$ .

With this aim in view, the following functions are introduced: S(t, x(t)) describing a continuous replenishment of the current stock and C(t, x(t)) describing a continuous distribution of the current stock. Then the difference S(t, x(t)) - C(t, x(t)) is a measure of the change of the current stock volume, i.e.

$$\frac{dx(t)}{dt} = S(t, x(t)) - C(t, x(t)).$$
<sup>(2)</sup>

Let us work out the functions that make up the right side of the equation (2), namely functions S(t, x(t)) and C(t, x(t)) in detail. The function of continuous replenishment S(t, x(t)) consists of three additive components, namely, from regulated replenishment of the stock, which is designated as  $S_{reg.}(\cdot)$ ; from unregulated replenishment  $S_{unreg.}(\cdot)$ ; and from random replenishment (for instance, a random stock replenishment due to an exceptionally high quality of productions or because of an expected sudden deficit of particular productions, etc.), which can be described mathematically as a random quantity  $X_s(t)$  that designating the total volume of productions that have been delivered into a particular warehouse from random and/or non-random sources by the time t. It is assumed for all types of replenishment that all orders are instantaneously executed, i.e. the shipping time for particular supplies is not considered in the present work. Let us interpret the introduced functions:

1) the function  $S_{reg.}(\cdot)$  can be interpreted as "one hundred per cent" (guaranteed) constant replenishment of the current stock of divisible productions, i.e. replenishment of the current stock that takes place regularly according to a contract during the segment  $[T_s, T_e]$ , with the volume of such replenishment being either constant (i.e.  $S_{reg.} \equiv const.$ ) or depending on t (i.e. being a function of the argument time  $S_{reg.} = S_{reg.}(t)$ , or else being functionally dependent on x(t) (i.e.  $S_{reg.} = S_{reg.}(t, x(t))$ );

2) the function  $S_{unreg.}(\bullet)$  obviously depends on t and functionally on x(t), and also on a certain quantity  $x_0(t)$ , which designates the minimal volume of stock in a particular warehouse necessary for administering unregulated stock replenishment on condition that such replenishment is guaranteed. In other words,  $S_{unreg.} = S_{unreg.}(t, x(t), x_0(t)) = k_0 \cdot x(t) \cdot \delta(x(t), x_0(t))$ , where  $k_0$  is a proportion coefficient, and the function  $\delta(x(t), x_0(t))$  is an indicator function, which has the form

$$\delta(x(t), x_0(t)) = \begin{cases} 1, & \text{if } x(t) \le x_0(t), \\ 0, & \text{if } x(t) > x_0(t); \end{cases}$$
(3)

3) the random quantity  $X_s(t)$  determines the total volume of productions that was delivered into the warehouse by the time t due to random circumstances from random and/or non-random sources. Then the quantity  $X_s(t + \Delta t)$  designates the sum total of all random deliveries by the time t + dt, where dt is an elementary interval of time (on analogy with the terminology of mathematical physics), and  $0 < dt \ll 1$ ,  $t + dt \le T_e$ . Consequently, it is possible to introduce a stochastic differential of a random process  $X_s(t)$ , namely, the quantity

$$dX_{s}dt \stackrel{def}{=} X_{s}(t+dt) - X_{s}(t),$$

which determines a random addition to the current stock of divisible productions during the elementary interval of time dt.

Now the function C(t, x(t)) that is contained in the right-hand side of the equation (2) and describes the dynamics of the continuous distribution of the current stock of divisible productions can be looked at in more detail. The function of continuous distribution C(t, x(t)) consists of three additive components, regulated distribution which is marked as  $C_{reg.}(\cdot)$ ; unregulated distribution  $C_{unreg.}(\cdot)$ , and random distributions (similar to random replenishment, there can be circumstances due to which random distribution takes place) that can be mathematically presented as a random quantity  $X_C(t)$  designating the total volume of productions that was taken away from the warehouse by the time t due to random circumstances. Let us now interpret the introduced functions.

1) the function  $C_{reg.}(\cdot)$  can be interpreted as "strong" (guaranteed) constant distribution of the current stock of divisible productions, i.e. the volume of the current stock that is regularly taken away from the warehouse according to contracts during the segment  $[T_s, T_e]$ , with the volume of such distribution being either constant (i.e.  $C_{reg.} \equiv const.$ ) or depending on t (i.e. being a function of the argument time  $C_{reg.} = C_{reg.}(t)$ , or else being functionally dependent on x(t) (i.e.  $C_{reg.} = C_{reg.}(t, x(t))$ );

2) the function  $C_{unreg.}(\cdot)$  depends on the time t and functionally on x(t) in general, as well as on a certain threshold function  $x_1(t)$ , which determines the stock volume of divisible productions allowing for its unregulated distribution,  $C_{unreg.} = C_{unreg.}(t, x(t), x_1(t))$ . In order to find an analytical expression of the function  $C_{unreg.}(t, x(t), x_1(t))$  the following assumptions can be made:

under  $x(t) \to \infty$  must be  $C_{unreg.}(t, x(t), x_1(t)) \to k_1$ , where the quantity  $k_1$  is the capacity of distributing the stock volume of divisible productions from the warehouse in the sense that whatever the stock

replenishment (i.e. the quantity  $\max_{t \in [T_s, T_e]} S(t, x(t))$ , the warehouse can not possibly distribute the stock of divisible productions measured as  $k_1$  during the entire considered time segment  $[T_s, T_e]$ ; under  $x(t) \rightarrow k_2 \equiv const$ , where  $k_2$  is an averaged value of the replenishment volume that allows for unregulated distribution, must be

$$C_{unreg.}(t, x(t), x_1(t)) \to \begin{cases} 0, & \text{if } x_1(t) \ge k_1, \\ \frac{k_1}{2}, & \text{if } x_1(t) < k_1. \end{cases}$$

The last two suppositions allow for determining unknown analytical form of the function  $C_{unreg.}(t, x(t), x_1(t))$ :

$$C_{unreg.}(t, x(t), x_1(t)) = k_1 \cdot x(t) \cdot \frac{1 - \delta(x(t), x_1(t))}{x(t) + k_2},$$

where the indicator function  $\delta(x(t), x_1(t))$  has the same sense/value as in determining the function  $S_{unreg.}(t, x(t), x_0(t))$ ; it is derive by formula (3) with the corresponding substitution of  $x_1(t)$  for  $x_0(t)$ ;

3) the random quantity  $X_c(t)$  designates the total volume of productions that has been removed from the warehouse by the time t due to random circumstances. Then  $X_c(t + \Delta t)$  designates the sum total of random distribution by the time t + dt, where dt is an elementary interval of time, with  $0 < dt \ll 1$ ,  $t + dt \le T_e$ . It follows that a stochastic differential of a random process  $X_c(t)$  can be introduced, namely the quantity  $dX_c dt \stackrel{def}{=} X_c(t + dt) - X_c(t)$ , which designates a random distribution of the current stock of divisible productions during the elementary interval of time dt.

Thus, taking into account the above specification of functions S(t, x(t)) and C(t, x(t)) the differential the equation (2) takes on form

$$dx(t) = S_{reg.}(t, x(t))dt + k_0 \cdot x(t) \cdot \delta(x(t), x_0(t))dt + dX_s - -C_{reg.}(t)dt - k_1 \cdot x(t) \cdot \frac{1 - \delta(x(t), x_1(t))}{x(t) + k_2}dt - dX_c.$$
(4)

The following initial condition (5) must be added to (4):

$$\left. x(t) \right|_{t=T_s} = x_s. \tag{5}$$

The obtained equation (4) is the stochastic differential equation with respect to unknown random volume x(t) of the current stock of divisible productions; and this equation together with the initial condition (5) constitutes the Cauchy problem for determine required volume x(t) of the current stock of divisible productions.

It is significant that the summands  $dX_s$  and  $dX_c$  in the right-hand side of the equation (4) are not differentials in the usual sense; these summands must be understood in the sense of the Ito stochastic differential (see [7]). Besides, the indicator functions  $\delta(x(t), x_0(t))$  and  $\delta(x(t), x_1(t))$  in the righthand side of the equation (4), derived according to formula (3), are not differentiated functions, which is caused by non-differentiability of the functions S(t, x(t)) and C(t, x(t)). Consequently, the requirement  $\exists x'(t) \ \forall t \in [T_s, T_e]$ , which was identified in the beginning of this section as a necessary condition for mathematical correctness of the model, will not be met. That is why in order to render a mathematical sense to the stochastic differential equation (4), it is necessary to introduce into is a corresponding amendment-condition. An easily realizable amendment might be substitution of the scalar functions  $\delta(x(t), x_0(t))$  and  $\delta(x(t), x_1(t))$  by the corresponding quadratic functions (which are smooth functions)

on the sections  $[0, x_0(t)]$  and  $[0, x_1(t)]$ , respectively. Such substitution is easily performed on the ground of natural and apparent requirements

$$\hat{\delta}(x(t), x_0(t))\Big|_{x(t)=0} = 1; \ \hat{\delta}(x(t), x_1(t))\Big|_{x(t)=0} = 1; \ \hat{\delta}(x(t), x_0(t))\Big|_{x(t)=x_0(t)} = 0; \ \hat{\delta}(x(t), x_1(t))\Big|_{x(t)=x_1(t)} = 0;$$

and in the result the following differential functions are obtained:

$$\hat{\delta}(x(t), x_0(t)) = -\frac{3}{x_0^2(t)} \cdot x^2(t) + \frac{2}{x_0(t)} \cdot x(t) + 1 \text{ when } x(t) \in [0, x_0(t)],$$
  
$$\hat{\delta}(x(t), x_1(t)) = -\frac{3}{x_1^2(t)} \cdot x^2(t) + \frac{2}{x_1(t)} \cdot x(t) + 1 \text{ when } x(t) \in [0, x_1(t)].$$

It is obvious that other substitutions-approximations are possible (for instance, by splines, etc.), which in comparison to the described above approach, i.e. approximation of scalar functions  $\delta(x(t), x_0(t))$  and  $\delta(x(t), x_1(t))$  by the corresponding smooth functions  $\hat{\delta}(x(t), x_0(t))$  and  $\hat{\delta}(x(t), x_1(t))$  provide a higher level of precision. In this sense, there is certain ambiguity in determining the functions  $\hat{\delta}(x(t), x_0(t))$  and  $\hat{\delta}(x(t), x_1(t))$ , and hence ambiguity of the right-hand side of the equation (4). Thus, instead of the differential equation (4) having no mathematical sense a mathematically correctly formulated differential equation can be written down:

$$dx(t) = S_{reg.}(t, x(t))dt + k_0 \cdot x(t) \cdot \delta(x(t), x_0(t))dt + dX_s$$
  
- $C_{reg.}(t)dt - k_1 \cdot x(t) \cdot \frac{1 - \delta(x(t), x_1(t))}{x(t) + k_2}dt - dX_c.$ 

It is important to note the following with regard to the obtained stochastic differential equation. It is obvious that stochastic differentials of the random processes  $X_s(t)$  and  $X_c(t)$  can be conjoined if a random quantity X(t) designating the total volume of productions that were delivered to and distributed from, the warehouse by the time t due to random circumstances. Then we can indeed determine a stochastic differential of the random process X(t) as

$$dXdt \stackrel{def}{=} X(t+dt) - X(t),$$

and this quantity will determine the change dynamics of the random volume of the divisible productions' stock during the elementary interval of time dt, namely dXdt > 0 designates a random replenishment of stock during the elementary interval of time dt, and dXdt < 0 designates a random distribution of stock during the elementary interval of time dt. With this specification in taken into account, the last differential equation takes the following final form:

$$dx(t) = S_{reg.}(t, x(t))dt + k_0 \cdot x(t) \cdot \hat{\delta}(x(t), x_0(t))dt - -C_{reg.}(t)dt - k_1 \cdot x(t) \cdot \frac{1 - \hat{\delta}(x(t), x_1(t))}{x(t) + k_2}dt + dX,$$
(6)

where  $t \in [T_s, T_e]$ , functions  $S_{reg.}(t, x(t))$ ,  $C_{reg.}(t)$  and  $x_i(t)$  (i = 0, 1), as well as numerical parameters  $k_i$   $(i = \overline{0, 2})$  have the described above values and are viewed as the given initial data of the problem under consideration; the functions  $\hat{\delta}(x(t), x_0(t))$  and  $\hat{\delta}(x(t), x_1(t))$  are determined by the following formulas:

$$\hat{\delta}(x(t), x_0(t)) = \begin{cases} 0, & \text{if } x(t) > x_0(t), \\ -\frac{3}{x_0^2(t)} \cdot x^2(t) + \frac{2}{x_0(t)} \cdot x(t) + 1, & \text{if } x(t) \le x_0(t), \end{cases}$$
(7)

$$\hat{\delta}(x(t), x_{1}(t)) = \begin{cases} 0, & \text{if } x(t) > x_{1}(t), \\ -\frac{3}{x_{1}^{2}(t)} \cdot x^{2}(t) + \frac{2}{x_{1}(t)} \cdot x(t) + 1, & \text{if } x(t) \le x_{1}(t). \end{cases}$$
(8)

The stochastic differential equation (6) together with the initial condition (5), the initial given data  $S_{reg.}(t, x(t))$ ,  $C_{reg.}(t)$  and  $k_i$   $(i = \overline{0, 2})$ , as well as approximating smooth indicator functions (7) and (8) is the Cauchy stochastic problem. It is a stochastic mathematical model for determining the current stock volume of divisible homogeneous production. Unfortunately, the given paper did not investigate the issue of finding an analytical solution of the constructed model (5)–(8). Nevertheless, as the following section will demonstrate, if we additionally require that the random process X(t) will be the Markov random process, then the constructed continuous model (5)–(8) can be easily realized numerically (see, for instance, [5]).

**Remark 1.** Stochastic equation (6) shows that irrespective of the sign of the quantity  $x_s = x(t)|_{t=T}$ 

(i.e. irrespective of the initial condition (5)), unknown function x(t) can assume a negative value, which,

at first sight, does not make any economic sense. But a possibility of such a case was purposefully taken into account prior to constructing mathematical model (5)–(8), and this case can be understood as a debt of the warehouse with regard to the current stock of divisible productions. Besides, a closer look at the left-hand side of the equation (6) (as well as the equations (2) and (4)), it becomes obvious that there can be a case when  $\frac{dx(t)}{dt} < 0$ , which means a negative rate if the quantity  $\frac{dx(t)}{dt}$  is treated as the speed of a material point of the continuous medium in metric space, which has no physical sense. But if the quantity  $\frac{dx(t)}{dt}$  in the considered problem designates the change rate of the volume x(t) of the current stock

at the time  $t \in [T_s, T_e]$ , then the case  $\frac{dx(t)}{dt} < 0$  corresponds to the situation whereby the volume x(t) as a function of the time argument is a decreasing function, i.e. the accumulated stock of divisible

#### 3. Construction of Finite-Differenced Model for Determination of Random Volume of Divisible Homogeneous Production

In this section we offer a finite-differenced approximation of the mathematical model (5)–(8) for determination of current stock volume of divisible homogeneous production, which was constructed in the previous section. Besides, given some assumptions, we put forward a recurrent implicit differenced scheme for numeric determination of the random volume of divisible homogeneous production at given discrete moments of time.

Let us introduce a determinate (i.e. non-random) function

productions in the warehouse is decreasing.

$$f(t, x(t)) \stackrel{\text{def}}{=} \begin{cases} S_{\text{reg.}}(t, x(t)), & \text{if } x(t) > x_0(t), \\ S_{\text{reg.}}(t, x(t)) - \frac{3 \cdot k_0}{x_0^2(t)} \cdot x^3(t) + \frac{2 \cdot k_0}{x_0(t)} \cdot x^2(t) + k_0 \cdot x(t), & \text{if } x(t) \le x_0(t) \end{cases} - \\ - \begin{cases} C_{\text{reg.}}(t, x(t)) + k_1 \cdot \frac{x(t)}{x(t) + k_2}, & \text{if } x(t) > x_1(t), \\ C_{\text{reg.}}(t, x(t)) + \frac{3 \cdot k_1}{x_1^2(t)} \cdot \frac{x^3(t)}{x(t) + k_2} - \frac{2 \cdot k_1}{x_0(t)} \cdot \frac{x^2(t)}{x(t) + k_2}, & \text{if } x(t) \le x_1(t). \end{cases}$$

$$(9)$$

Then the stochastic equation (6) can be rewritten in a more compact way:

$$dx(t) = f(t, x(t))dt + dX(t),$$
<sup>(10)</sup>

and this equation is a particular instantiation (namely,  $f_1(t, x(t)) \equiv f(t, x(t))$ ;  $f_2(t, x(t)) \equiv 1$ ) of a more general stochastic differential equation in the Ito form

$$dx(t) = f_1(t, x(t))dt + f_2(t, x(t))dX(t),$$
(11)

where the functions  $f_i(t, x(t))$  (i = 1, 2) are supposed to be non-random functions, the random process X(t) the Markov random process X(t), and the quantity dX(t) is understood in the sense of a stochastic differential Markov random process X(t).

Under the mentioned assumptions, the Ito stochastic differential equation (11) allows for the following interpretation: for the stochastic differential dX(t), which is contained in the right-hand side of the equation (11), the quantity X(t) can be understood as a realized random quantity which assumes the given value  $\tilde{x} = X(\tilde{t})$  at the moment  $\tilde{t} \in [T_s, T_e]$ . Moreover, due to the assumption that X(t) is the Markov process the random quantity  $X(\tilde{t} + dt) = \tilde{x}$ , where  $0 < dt \ll 1$ ,  $\tilde{t} + dt \leq T_e$ , has a density of probability  $\rho(\tilde{x}) = \rho(\tilde{t}, \tilde{x}; \tilde{t} + dt, \tilde{x})$ . Then, if randomness of  $\tilde{t}$  is taken into account, the above speculation holds for  $\forall t \in [T_s, T_e]$ , i.e. for random  $t \in [T_s, T_e]$ , the random quantity  $X(t+dt) = \tilde{x}$  where  $0 < dt \ll 1$ ,  $t + dt \leq T_e$ , is determined by the density of probabilities  $\rho(\tilde{x}) = \rho(t, \tilde{x}; t + dt, \tilde{x})$  only if the random quantity X(t) assumed the concrete value  $\tilde{x}$  at the moment  $t \in [T_s, T_e]$ , i.e. if  $X(t) = \tilde{x}$ . This interpretation of the Ito stochastic differential equation (11) allows for rewriting the equation (11) in the finite-difference approximation, namely

$$\begin{aligned} x(t+\Delta t) - x(t) &= f_1(t, x(t)) \cdot \Delta t + f_2(t, x(t)) \cdot \left(X(t+\Delta t) - X(t)\right) = \\ &= f_1(t, x(t)) \cdot \Delta t + f_2(t, x(t)) \cdot \left(X(t+\Delta t) - \tilde{x}\right). \\ &\text{If we accept} \end{aligned}$$

$$x(t) = f_2(t, x(t)) \cdot X(t) = f_2(t, x(t)) \cdot \tilde{x}$$

then we obtain a recurrent correlation

$$x(t + \Delta t) = f_1(t, x(t)) \cdot \Delta t + f_2(t, x(t)) \cdot X(t + \Delta t),$$

which can be used for a discrete definition of the value of unknown function x(t). Indeed, if we break down the time segment  $[T_s, T_e]$  into N elementary time spaces of the length  $\Delta t_i$   $(i = \overline{0, N-1})$  we will obtain the discrete mesh

$$\hat{T} \stackrel{\text{def}}{=} \left\{ t_i : t_{i+1} = t_i + \Delta t_i \ \left( i = \overline{0, N-1} \right), \ t_0 = T_s, \ t_N = T_e \right\}$$

and after designating

$$x_i \stackrel{def}{=} x(t_i), \quad \tilde{x}_i \stackrel{def}{=} X(t_i) = \frac{x(t_i)}{f_2(t_i, x(t_i))}$$

it is possible to write down the following recurrent implicit differenced scheme for determining the quantity x(t) numerically:

$$x_{i+1} = f_1(t_i, x_i) \cdot \Delta t_i + f_2(t_i, x_i) \cdot \tilde{x}_{i+1},$$
(12)

where random quantities  $\tilde{x}_{i+1}$  are determined by the density of probabilities  $\rho\left(t_i, \frac{x_i}{f_2(t_i, x_i)}; t_{i+1}, \tilde{\tilde{x}}\right)$ .

**Remark 2.** Mathematical model (5)–(8) constructed in the Section 2 can be solved analytically with the help of integrals of Stratonovich and Ito (see [7]) assuming a Markov nature of the random process X(t). If this

assumption is not made (or can not be made due to specificity of the particular task of inventory control), the question of how to analytically integrate the stochastic differential equation (8) remains, unfortunately, still open, and as mentioned before research into this issue was not undertaken in the present paper. As shown in [8], though, with certain additional conditions but without assuming the Markov nature of the random process X(t) an effective approximation of a stochastic differential equation such as (11), particularly the equation (10), which is the equation (6) in the mathematical model (5)–(8) constructed in the Section 2.

**Remark 3.** The constructed recurrent differenced scheme (12) together with the initial condition (5) is a finite differenced mathematical model for defining one of possible trajectories of the random quantity x(t), i.e. the constructed finite differenced model (12); (5) allows for defining approximate

values of the quantity x(t) at the moments of time  $t_i$ :  $t_{i+1} = t_i + \Delta t_i$   $(i = \overline{0, N-1})$ ,  $t_0 = T_s$ ,  $t_N = T$ .

#### 4. Stochastic Continuous Models for Defining Volumes of Current Stock of Divisible Productions at Several Interconnected Warehouses Simultaneously

The present section suggests two stochastic continuous mathematical models for defining volumes of current stock of divisible homogeneous and heterogeneous production at several interconnected warehouses simultaneously. For achieving this aim, similarly to the Section 2, apparatus of mathematical physics is used and principle of continuous medium, the language of the theory of partial differential equations is chosen as a modelling language. Because of paper's space limitations there is, unfortunately, no opportunity to present the entire chain of argumentation and all calculations related to constructing these models "from scratch"; they are only mathematically represented in what follows, with minimal explanation.

So,  $m \in \mathbb{N}$  warehouses are under consideration, and it is assumed that dynamics of the volume of divisible homogeneous production in all m warehouses is subject to the stochastic differential equation (6) which was obtained in the Section 2. For the stochastic differential dX(t) that is contained in the right-hand side of the equation (6), the quantity X(t) will be viewed as a realized random quantity which assumed the given value  $\tilde{x} = X(t)$  at the time  $t \in [T_s, T_e]$ , i.e. the given warehouse has the volume of divisible homogeneous production  $\tilde{x} = x(t)$  at the fixed time  $t \in [T_s, T_e]$ ; as the course of constructing equation (6) shows, this volume can comprise both determined and random constituent volumes. Then the random quantity  $X(t+dt) = \tilde{\tilde{x}}$ , where  $0 < dt \ll 1$ ,  $t+dt \leq T_e$ , designates a random volume of homogeneous production was present in this very warehouse at the previous moment t. Consequently, it can be said that the random quantity X(t+dt) has the density of probability  $\rho(\tilde{\tilde{x}}) = \rho(t, \tilde{x}; t+dt, \tilde{\tilde{x}})$ .

Since the continuous mathematical model (5)–(8) constructed in the Section 2 assumed the existence of one warehouse where there was volume  $x_s = x(t)\Big|_{t=T_s}$ , of divisible homogeneous foods at the initial moment of time  $t = T_s$ , for *m* interconnected warehouses there are obviously *m* initial conditions  $x^{\{i\}}(t)\Big|_{t=T_s} = x_s^{\{i\}}, \ i = \overline{1,m}$ , where  $x^{\{i\}}(t)$  designates a random volume of the divisible homogeneous production in an *i*-warehouse at the time  $t \in [T_s, T_e]$ . That is why once these random initial volumes  $x_s^{\{i\}}, \ i = \overline{1,m}$  were distributed on the axis *OX* of the Cartesian rectangular system of coordinates, these irregularly distributed initial volumes can be mentally identified with the distribution of the warehouses on the axis *OX*. This identification allows for constructing the required mathematical model. It is worth mentioning here that topology of the imagined distribution of warehouses on the axis *OX* does not have to match the typology of distributing initial quantities-volumes  $x_s^{\{i\}}$ ,  $i = \overline{1,m}$ ; this is natural and obvious.

After the above mentioned identification we have a certain set of interconnected warehouses (SIW), and we can construct a mathematical model for establishing the dynamics of random volumes of divisible homogeneous production in this SIW ignoring the dynamics of a random volume of divisible homogeneous production in any individual warehouse.

Let us consider a relatively short segment  $[x(t), x(t) + \Delta x(t)]$  of the length  $\Delta x(t)$  and introduce the functional  $\Delta \Psi(t, x(t))$  of the function- volume x(t) which describes the number of elements SIW that can be found in the segment  $[x(t), x(t) + \Delta x(t)]$ . In other words,  $\Delta \Psi(t, x(t))$  is the number of warehouses distributed on a short segment  $[x(t), x(t) + \Delta x(t)]$  of the length  $\Delta x(t)$ .

Then  $\frac{\Delta \Psi(t, x(t))}{\Delta x(t)}$  can be treated as probability of the warehouse with the volume x(t) of productions

being on the segment  $[x(t), x(t) + \Delta x(t)]$ . Consequently, we can move over to the limit with  $\Delta x(t) \rightarrow 0$  and define a new function

$$p(t,x(t)) \stackrel{\text{def}}{\equiv} \lim_{\Delta x(t) \to 0} \frac{\Delta \Psi(t,x(t))}{\Delta x(t)},$$

which is the density of distribution of warehouses according to random volumes x(t) of divisible homogeneous production. Then the function

$$\Psi(t) \stackrel{\text{def}}{=} \int_{x_1}^{x_2} p(t, x(t)) dx(t)$$

designates the number of warehouses with random volumes  $x(t) \in [x_1(t), x_2(t)]$  at the time moment  $t \in [T_s, T_e]$ .

It is easily seen that

$$\int_{T_s}^{T_e} \Psi(t) dt \equiv m; \quad \int_{-\infty}^{+\infty} p(t, x(t)) dx(t) \equiv 1.$$

Now the density of distribution p(t, x(t)) of warehouses according to random volumes x(t) of divisible homogeneous production is defined, and we can establish the law of distributing warehouses according to random volumes, i.e. to find out the rule that governs the change of the function p(t, x(t)). For this, the axis OX is divided into two parts, an arbitrary segment  $[x_1(t), x_2(t)]$  and the view of this segment, i.e. the domain  $(-\infty, x_1(t)) \cup (x_2(t), +\infty)$ . As random volumes of productions in warehouses change with the course of time, it will mean in our case that warehouses will be moving along the axis OX in this course of time. This, in turn, means that during the segment of time  $[t_1, t_2]$ ,  $\forall t_1, t_2 \in [T_s, T_e]$  a certain number of warehouses will have random volumes of divisible homogeneous productions that are no bigger than  $x_1(t)$  and no less than  $x_2(t)$ , i.e. some warehouses will be located in the segment  $[x_1(t), x_2(t)]$  whereas their remaining number will be outside this segment, or in the domain  $(-\infty, x_1(t)) \cup (x_2(t), +\infty)$ . Thus it will be quite correct if the equation of balance of warehouses for the segment  $[x_1(t), x_2(t)]$  in the segment of time  $[t_1, t_2]$  is presented in the following way (on analogy with a widely known approach in mathematical physics whereby mathematical models are constructed for heat conductivity, waves, diffusion, radiation, and other physical processes):

$$\Delta \Psi(t_1, t_2) \stackrel{\text{def}}{=} \Psi(t_2) - \Psi(t_1) = \Psi^{\{1\}} + \Psi^{\{2\}}, \tag{13}$$

where  $\Psi^{\{1\}}$  is the number of warehouses located in the segment  $[x_1(t), x_2(t)]$  in the segment of time  $[t_1, t_2]$  due to non-random replenishments and distributions of divisible homogeneous productions;  $\Psi^{\{2\}}$  is the number of warehouses located in the segment  $[x_1(t), x_2(t)]$  in the segment of time  $[t_1, t_2]$  due to random replenishments and distributions of divisible homogeneous production; and the function  $\Delta\Psi(t_1, t_2)$  in the left-hand side of (16) is calculated according to the formula

$$\Delta \Psi(t_1, t_2) \stackrel{\text{def}}{=} \Psi(t_2) - \Psi(t_1) = \int_{x_1}^{x_2} p(t_2, x) dx - \int_{x_1}^{x_2} p(t_1, x) dx =$$

$$= \int_{x_1}^{x_2} p(t, x(t)) \Big|_{t=t_1}^{t=t_2} dx = \int_{x_1}^{x_2} dx \int_{t_1}^{t_2} \frac{\partial p(t, x(t))}{\partial t} dt.$$
(14)

It is obvious that the quantities  $\Psi^{\{i\}}$  (i = 1, 2) can be negative, and this is then treated as a removal of warehouses from the segment  $[x_1(t), x_2(t)]$ . The final formulas for the functions  $\Psi^{\{i\}}$  and  $\Psi^{\{2\}}$  are given below without conclusion (there is an elegant conclusion which is not given here due to the space constraints):

$$\Psi^{\{1\}} = \int_{t_1}^{t_2} \left\{ p\left(t, x_1\left(t\right)\right) \cdot \vartheta\left(t, x_1\left(t\right)\right) - p\left(t, x_2\left(t\right)\right) \cdot \vartheta\left(t, x_2\left(t\right)\right) \right\} dt = \int_{t_1}^{t_2} dt \int_{x_1}^{x_2} \frac{\partial}{\partial x(t)} \left( p\left(t, x(t)\right) \cdot \vartheta\left(t, x(t)\right) \right) dx(t),$$
(15)

$$\Psi^{\{2\}} = \int_{t_1}^{t_2} dt \int_{x_1}^{x_2} \left\{ -\frac{\partial}{\partial x} \left( a\left(x(t), t\right) \cdot p\left(t, x(t)\right) \right) + \frac{1}{2} \cdot \frac{\partial^2}{\partial x^2} \left( b\left(x(t), t\right) \cdot p\left(t, x(t)\right) \right) \right\} dx(t),$$
(16)

where the function  $\rho(z,s;x,t)$  is a transitional function of the probability density of a diffusion stochastic process X(t) (for instance, see [9], [10]); the function  $\mathcal{G}(t,x(t))$  designates the change rate of the random volume x(t) of the current stock of divisible homogeneous production in the set of interconnected warehouses (SIW) at the time t, and is determined by the stochastic equation

$$\mathscr{G}(t,x(t)) = \frac{dx(t)}{dt} = S(t,x(t)) - C(t,x(t)),$$

where the functions S(t, x(t)) and C(t, x(t)) have the same values as mentioned in the Section 2. The functions a(x(t), t) and b(x(t), t) are calculated by the formulas

$$a(t,x(t)) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \cdot \int_{|x(t)-z(t)| \leq \varepsilon} (x(t)-z(t)) \cdot \rho(x(t),t;z(t),t+\Delta t) dz(t) \quad (\forall \varepsilon > 0, \forall z \in \mathbb{R}^{1}),$$
  
$$b(t,x(t)) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \cdot \int_{|x(t)-z(t)| \leq \varepsilon} (x(t)-z(t))^{2} \cdot \rho(x(t),t;z(t),t+\Delta t) dz(t) \quad (\forall \varepsilon > 0, \forall z \in \mathbb{R}^{1}).$$

Taking into account expressions (14)–(16) in formula (13), the following equation is obtained:

$$\int_{x_{1}}^{x_{2}} dx \int_{t_{1}}^{t_{2}} \frac{\partial p(t, x(t))}{\partial t} dt = \int_{t_{1}}^{t_{2}} dt \int_{x_{1}(t)}^{x_{2}} \frac{\partial}{\partial x(t)} \left( p(t, x(t)) \cdot \vartheta(t, x(t)) \right) dx(t) +$$

$$+ \int_{t_{1}}^{t_{2}} dt \int_{x_{1}}^{x_{2}} \left\{ -\frac{\partial}{\partial x} \left( a(t, x(t)) \cdot p(t, x(t)) \right) + \frac{1}{2} \cdot \frac{\partial^{2}}{\partial x^{2}} \left( b(t, x(t)) \cdot p(t, x(t)) \right) \right\} dx(t),$$

which due to arbitrariness of the selected volume segment  $[x_1(t), x_2(t)]$ , arbitrariness of the selected time segment  $[t_1, t_2]$ , and in accordance with the First Mean Value Theorem (use of this theorem here is quite rightful because all its requirements are met) can be written in the following way:

$$\frac{\partial p(t,x(t))}{\partial t} = -\frac{\partial}{\partial x} \left( \left[ a(t,x(t)) + \vartheta(t,x(t)) \right] \cdot p(t,x(t)) \right) + \frac{1}{2} \cdot \frac{\partial^2}{\partial x^2} \left( b(t,x(t)) \cdot p(t,x(t)) \right).$$
(17)

The resulting stochastic equation (17) is the parabolic type particular differential equation, and together with the above mentioned functions  $a(t, x(t)) \neq 0$ ,  $b(t, x(t)) \neq 0$  and  $\vartheta(t, x(t))$ , as well as corresponding initial and boundary conditions (for instance, the Newton type boundary conditions, or Neumann boundary conditions, or non-located boundary conditions) it makes the required mathematical model for determining unknown density of distribution p(t, x(t)) of exactly  $m \in \mathbb{N}$  warehouses according to random volumes x(t) of divisible homogeneous production. It is not difficult to see that the equation (17) is a particular case of the widely known Kolmogorov's equation for the Markov stochastic process X(t) with a transition function of the density of probability  $\rho(z, s; x, t)$ .

The next stochastic continuous model (with the Dirichlet boundary conditions) is an informal generalization (the corresponding conclusion is rather complex and therefore not presented in the given article) of the above mentioned model: it describes the dynamics of unknown density of distribution  $p(t, x_1(t), ..., x_n(t))$  of exactly  $m \in \mathbb{N}$  warehouses according to random volumes  $x(t) = (x_1(t), ..., x_n(t))$  of divisible  $n \in \mathbb{N}$  heterogeneous products

$$\begin{split} \frac{\partial p(t, x(t))}{\partial t} &= \frac{1}{2} \cdot \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{\partial^{2}}{\partial x_{i} \partial x_{j}} \left( b_{ij}\left(t, x(t)\right) \cdot p(t, x(t)) \right) - \sum_{i=1}^{n} \frac{\partial}{\partial x_{i}} \left( \left[ a_{i}\left(t, x(t)\right) + \mathcal{G}_{i}\left(t, x(t)\right) \right] \cdot p(t, x(t)) \right), \\ p(t, x(t)) \Big|_{t=T_{s}} &= p_{0}\left(x_{1}, \dots, x_{n}\right), \ x_{i} \in \mathbb{R}^{1} \quad \forall i = \overline{1, n}; \\ p(t, x(t)) \Big|_{x_{i}(t) = l_{i}^{[1]} - 0} &= p_{i}^{[1]}\left(t\right), \ l_{i}^{[1]} \in \mathbb{R}^{1} \quad \left(i = \overline{1, n}\right); \\ p(t, x(t)) \Big|_{x_{i}(t) = l_{i}^{[2]} - 0} &= p_{i}^{[2]}\left(t\right), \ l_{i}^{[2]} \in \mathbb{R}^{1} \quad \left(i = \overline{1, n}\right), \end{split}$$

where  $x(t) = (x_1(t), ..., x_n(t)) \in \bigcup_{i=1}^n [I_i^{\{1\}}, I_i^{\{2\}}]$ , the function  $x_i(t)(i = \overline{1, n})$  describes the random volume of *i*-th divisible product at the time moment  $t \in [T_s, T_e]$ ; the function  $\mathcal{G}_i(t, x(t))$  describes the change rate of the random volume  $x_i(t)$  of the current stock of *i*-th divisible product in the set  $m \in \mathbb{N}$  of interconnected warehouses at the time moment t; the functions  $a_i(x(t), t)(i = \overline{1, n})$  and  $b_{ij}(t, x(t))(i = \overline{1, n}; j = \overline{1, n})$  are calculated according to the formulas

$$a_{i}(t,x(t)) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \cdot \int_{B_{\varepsilon}(z(t))} (x_{i}(t) - z_{i}(t)) \cdot \rho(x,t;z,t+\Delta t) dz_{1}(t) \dots dz_{1}(t),$$
  
$$b_{ij}(t,x(t)) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \cdot \int_{B_{\varepsilon}(z(t))} (x_{i}(t) - z_{i}(t)) \cdot (x_{j}(t) - z_{j}(t)) \cdot dz_{1}(t) \dots dz_{1}(t),$$

where the function  $\rho(x,t;z,t+\Delta t) \stackrel{\text{def}}{=} \rho(x_1(t),...,x_n(t),t;z_1(t),...,z_n(t),t+\Delta t)$  is a transition function of the density of probabilities of the diffusion stochastic process  $X(t) \stackrel{\text{def}}{=} (X_1(t),...,X_n(t))$  (for example, see [9–10] and the lists of corresponding literature in them), and  $B_{\varepsilon}(z(t)) \stackrel{\text{def}}{=} \{x(t): ||x(t)-z(t)||_{\mathbb{R}^n} \le \varepsilon\}$  is the closed  $\varepsilon$ -neighbourhood of the point.

#### Conclusions

- In the given work the stochastic continuous mathematical models for determining the random volumes of current stock of both homogeneous and heterogeneous divisible production in one or several interconnected warehouses are constructed. The constructed models can be used for on-line monitoring of the dynamics of the productions random volumes.
- In future research the authors intend to investigate the questions of the constructed mathematical models solvability, finding the minimal sufficient conditions to ensure the uniqueness of their solutions, as well as the development of the analytical and stable numerical methods for finding the solutions of the constructed mathematical models.

#### References

- 1. Kopytov, E., Muravjov, A., Greenglaz, L. and G. Burakov. Investigation of two strategies in inventory control system with random parameters. In: *Proceedings of the XXI European Conference on Modelling and Simulation (ECMS-2007), Prague, Czech Republic.* Prague, 2007, pp. 566–571.
- 2. Ashmanov, S. A. *Mathematical Models and methods in Economics*. Moscow: Lomonosov MSU Press, 1980. (In Russian)
- 3. Nikaido, H. Convex structures and economic theory. New York-London: Academic Press, 1968.
- 4. Tikhonov, A. N. and A. A. Samarsky. *Equations of Mathematical Physics*. Moscow: Lomonosov MSU Press, 2004. (In Russian)
- 5. Ito, K. Selected papers. New York: Springer Press, 1987.
- 6. Milstein, G. N. *Numerical integration of stochastic differential equations*. New York-London: Kluwer Academic Publishers, 1995.
- 7. Kuznetsov, D. F. *Stochastic Differential equations: theory and practice of numerical solutions.* St. Petersburg: Polytechnic University Press, 2007. (In Russian)
- 8. Diend, I. On approximation of Ito stochastic equations. *Mathematical Transactions*, Vol. 181, No 6, 1990, pp. 743–750. (In Russian)
- 9. Samarsky, A. A. and A. P. Mikhailov. Mathematical Modelling. Moscow: Fizmatlit Press. 2002.
- 10. Gikhman, I. I. and A. V. Skorokhod. *Stochastic differential equations and their applications*. Kiev: Naukova Dumka Publishers, 1982. (In Russian)

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 31–39 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# STATISTICAL INFERENCE USING ENTROPY BASED EMPIRICAL LIKELIHOOD STATISTICS

G. Gurevich<sup>1\*</sup>, A. Vexler<sup>2</sup>

<sup>1</sup>The Department of Industrial Engineering and Management, SCE- Shamoon College of Engineering Beer-Sheva 84100, Israel E-mail: gregoryg@sce.ac.il <sup>2</sup>Department of Biostatistics, The State University of New York at Buffalo Buffalo, NY 14214, USA E-mail: avexler@buffalo.edu

In this article, we show that well known entropy-based tests are a product of empirical likelihood ratio. This approach yields stable definitions of entropy-based statistics for goodness-of fit tests and provides a simple development of two-sample tests based on samples entropy that have not been presented in the literature. We introduce the distribution-free density-based likelihood techniques, applied to test for goodness-of-fit. In addition, we propose and examine nonparametric two-sample likelihood ratio tests for the case-control study based on samples entropy. The Monte Carlo simulation study indicates that the proposed tests compare favourably with the standard procedures, for a wide range of null/alternative distributions.

Keywords: empirical likelihood, entropy, goodness-of-fit tests, two-sample nonparametric tests, case-control study

#### 1. Introduction

The likelihood approach is a powerful and widely-used tool for parametric statistical inference. As an example, consider the simple hypothesis testing problem where given a sample of k independent identically distributed observations  $X_1, ..., X_k$ , we want to test the hypothesis

$$H_0: X_1, ..., X_k \sim F_0$$
 versus  $H_1: X_1, ..., X_k \sim F_1$ , (1)

where  $F_0$  and  $F_1$  are some distributions with density functions  $f_0(x)$  and  $f_1(x)$ , respectively. By virtue of the Neyman–Pearson Lemma, the most powerful test-statistic for (1) is the likelihood ratio

$$\prod_{i=1}^{k} f_1(X_i) \\
\prod_{i=1}^{k} f_0(X_i),$$
(2)

where density functions  $f_0(x)$  and  $f_1(x)$  are assumed to be completely known. However, if the alternative distribution  $F_1$  is not known, the hypotheses (1) define a goodness-of-fit problem. For this situation, to use the likelihood ratio statistic one need to estimate a likelihood function in numerator of (2). There has been much recent development of various empirical likelihood type approximations to parametric likelihood functions. The empirical likelihood (EL) method based on empirical distributions has been dealt with extensively in the literature (e.g., Owen [7]). The EL function has the form of  $L_p = \prod_{i=1}^{k} p_i$ , where the components  $p_i$ , i = 1, ..., k maximize the likelihood  $L_p$ , satisfying empirical constraints

<sup>\*</sup> Address for correspondence: Gregory Gurevich, The Department of Industrial Engineering and Management, SCE – Shamoon College of Engineering, Beer-Sheva 84100, Israel; e-mail: gregoryg@sce.ac.il.

(e.g.,  $\sum_{i=1}^{k} p_i = 1$  and  $\sum_{i=1}^{k} p_i X_i = 0$ ). Computation of  $p_i$ , i = 1,...,k is based on a simple exercise in Lagrange multipliers (for details, see Owen [7]). This nonparametric approach is a result of consideration of the 'distribution functions'-based likelihood  $\prod_{i=1}^{k} (F(X_i) - F(X_i -))$  over all distribution functions F. Taking into account that the Neyman–Pearson Lemma operates under 'density functions'-based forms of likelihood functions, Vexler and Gurevich [11] applied the main idea of the EL technique to construct density-based empirical estimation of the parametric likelihood  $L_f = \prod_{i=1}^{k} f(X_i)$ , where f(x) is a density function. They considered the likelihood function  $L_f$  in the form of  $L_f = \prod_{i=1}^{k} f(X_i) = \prod_{i=1}^{k} f(X_{(i)}) = \prod_{i=1}^{k} f_i$ , where  $f_i = f(X_{(i)})$ , and  $X_{(1)} \leq X_{(2)} \leq \ldots \leq X_{(k)}$  are the order statistics derived from  $X_1, \ldots, X_k$ . The estimators of  $f_i$ ,  $i = 1, \ldots, k$  that maximize  $L_f$  and satisfy some empirical constrains have the following form:  $f_i = 2m/(k(X_{(i+m)} - X_{(i-m)}))$ ,  $i = 1, \ldots, k$ . Therefore, the maximum EL method applied to (2) with known  $f_0(x)$  and unknown  $f_1(x)$  forms the test-statistic

$$T_{mk} = \frac{\prod_{i=1}^{k} \frac{2m}{k(X_{(i+m)} - X_{(i-m)})}}{\prod_{i=1}^{k} f_0(X_i)}.$$
(3)

Note

that

$$\log\left(\prod_{i=1}^{k} \frac{2m}{k} \left( \frac{X_{(i+m)} - X_{(i-m)}}{k} \right) \right) = -kH(m,k), \quad \text{where}$$

 $H(m,k) = k^{-1} \sum_{i=1}^{k} \log(k(X_{(i+m)} - X_{(i-m)})/2m) \text{ was presented by Vasicek [10], as an estimator of the entropy of the density <math>f(x)$ , for some m < k/2, i.e., the statistic H(m,k) estimates  $H(f) = E(-\log(f(X_1))) = -\int_{-\infty}^{+\infty} f(x)\log(f(x))dx = \int_{0}^{1} \log(\frac{d}{dp}F^{-1}(p))dp$ . The power of the tests based on the statistic T, strongly depends on values of m and this restricts applicability of (3)-type

on the statistic  $T_{mk}$  strongly depends on values of m and this restricts applicability of (3)-type test-statistics to real-data problems. Dealing with this problem and reconsidering their empirical constraints for density functions, Vexler and Gurevich [11] proposed the statistic  $T_k^* = \min_{1 \le m < k^{1-\delta}} \left( \prod_{i=1}^k \frac{2m}{k(X_{(i+m)} - X_{(i-m)})} / \prod_{i=1}^k f_0(X_i) \right)$  as a modification of the entropy based statistic  $T_{mk}$ .

Considering the problem (1) where, under the alternative hypothesis,  $f_1(x)$  is completely unknown, whereas, under the null hypothesis,  $f_0(x) = f_0(x; \theta)$  is known up to the vector of parameters  $\theta = (\theta_1, ..., \theta_d)$  (here,  $d \ge 1$  defines a dimension of the vector  $\theta$ ), they proposed the statistic

$$G_{k} = \min_{1 \le m < k^{1-\delta}} \frac{\prod_{i=1}^{k} \frac{2m}{k(X_{(i+m)} - X_{(i-m)})}}{\prod_{i=1}^{k} f_{0}(X_{i}, \hat{\boldsymbol{\theta}})},$$
(4)

where  $0 < \delta < 1$  and  $\hat{\theta}$  estimates  $\theta$  (e.g.,  $\hat{\theta}$  is the maximum likelihood estimator of  $\theta$ ). They also proved that if some general conditions are satisfied for density functions  $f_0(x)$ ,  $f_1(x)$  and

for the estimator  $\hat{\theta}$ , then under  $H_0$ ,  $k^{-1}\log(G_k) \xrightarrow{P} 0$ , while, under  $H_1$ ,  $k^{-1}\log(G_k) \xrightarrow{P} E\log\left(\frac{f_1(X_1)}{f_0(X_1;\mathbf{a})}\right) > 0$ , where  $\mathbf{a} = (a_1,...,a_d)$  is a vector with finite components, as  $k \to \infty$ . That is, with a test-threshold *C* related to the type I error  $\alpha = \sup_{\alpha} P_{H_0}(\log(G_k) > C)$  in mind,  $P_{H_1}(\log(G_k) > C) \longrightarrow 1$ . That means that a test based on the statistic  $G_k$  has the asymptotic power one, i.e., is a consistent test.

#### 2. Empirical Likelihood Ratio Tests for Uniformity and Normality

#### 2.1. Test for Uniformity

Consider a problem (1), where  $F_0 = Unif(0,1)$ ,  $F_1$  is an unknown distribution with a finite variance and continuous density function  $f_1(x)$  concentrated on the interval [0,1]. In accordance with (4), the suggested test is: reject  $H_0$  if

$$\min_{1 \le m < k^{1-\delta}} \prod_{i=1}^{k} \frac{2m}{k(X_{(i+m)} - X_{(i-m)})} > C , \qquad (5)$$

where  $0 < \delta < 1$ , *C* is a test-threshold.

Note that the statistic  $U_{mk} = \prod_{i=1}^{k} \frac{2m}{k(X_{(i+m)} - X_{(i-m)})}$  with a fixed m < n/2 was considered by

Dudewicz and van der Meulen [4] as a test-statistic of the entropy-based test for uniformity. This test is a very efficient decision rule provided that optimal values of m, subject to  $f_1(x)$  and k, are applied to the statistic  $U_{mk}$  (Dudewicz and van der Meulen [4]). In practice, since  $f_1(x)$  is completely unknown, we risk choosing m that leads to a  $U_{mk}$ -based test having the power that is lower than that of other known tests for uniformity (e.g., Zhang [12]). In contrast to this, a Monte Carlo study, presented by Vexler and Gurevich [11], demonstrates that, in many cases, the test (5) provides the power that is close to the power of  $U_{mk}$  -based tests with optimal m 's, calculated empirically.

Test-threshold C for the test (5) can be obtained exactly or approximately by simulations from

the equation 
$$P_{X_1,\ldots,X_k\sim Unif(0,1)}\left\{\log\left(\min_{1\leq m< k^{1-\delta}}\prod_{i=1}^k \left(\frac{2m}{k}\left(X_{(i+m)} - X_{(i-m)}\right)\right)\right) > C\right\} = \alpha$$
, for each desired significance level  $\alpha$ 

In accordance with the asymptotic properties of the statistic (4), presented in Section 1, the test (5) is consistent as  $k \to \infty$ .

#### 2.2. Test for Composite Hypothesis of Normality

Consider a problem (1), where  $F_0$  is a normal distribution with unknown expectation  $\mu$  and variance  $\sigma^2$ ,  $F_0 = Norm(\mu, \sigma^2)$ ,  $F_1$  is an unknown distribution with a finite variance and continuous density function  $f_1(x)$ . In accordance with (4), the suggested test is: reject  $H_0$  if

$$\min_{1 \le m < k^{1-\delta}} \left( 2\pi e s^2 \right)^{k/2} \prod_{i=1}^k \frac{2m}{k \left( X_{(i+m)} - X_{(i-m)} \right)} > C , \qquad (6)$$

where  $0 < \delta < 1$ ,  $s^2 = \frac{1}{k} \sum_{j=1}^{k} \left( X_j - \frac{1}{k} \sum_{k=1}^{k} X_k \right)^2$ , *C* is a test-threshold.

Note that the statistic 
$$N_{mk} = (2\pi e s^2)^{k/2} \prod_{i=1}^k (2m/k(X_{(i+m)} - X_{(i-m)}))$$
 is known, for  $m < n/2$ , to

be an efficient test statistic based on sample entropy (e.g., Vasicek [10], Arizono and Ohta [1]; Park and Park [8]). The tests for normality based on sample entropy are exponential rate optimal procedures (Tusnady [9]). This agrees with the fact that commonly likelihood ratio tests have optimal statistical properties and likelihood ratio type decision rules are simple in applications. The power of the test based on statistic  $N_{mk}$  strongly depends on values of m. Assuming information regarding the distribution functions of the alternative hypothesis, Monte Carlo simulation results, published in the relevant literature, point out values of m (subject to k) that provide high levels of the power of the test based on  $N_{mk}$ . However, since  $f_1(x)$  is completely unknown, we risk choosing m that leads to a  $N_{mk}$ -based test having the power that is lower than that of other known tests for normality (e.g., Vexler and Gurevich [11]). In this sense, the main advantage of the proposed test (6) for normality is that his statistic is not depend on unknown parameters (following Vexler and Gurevich [11]), we recommend the value of  $\delta = 0.5$  in definition of (6)).

Since, under  $H_0$ , the statistic of the proposed test (6) does not depend on values of  $\mu$  and  $\sigma^2$ , the test-threshold C for this test can be obtained exactly or approximately by simulations from the equation  $P_{X_1,...,X_k \sim Norm(0,1)} \left\{ log \left( \min_{1 \le m < k^{1-\delta}} (2\pi e s^2)^{k/2} \prod_{i=1}^k (2m/k (X_{(i+m)} - X_{(i-m)})) \right) > C \right\} = \alpha$ , for each

desired significance level  $\alpha$ .

In accordance with the asymptotic properties of the statistic (4), presented in Section 1, the test (6) is consistent as  $k \to \infty$ .

#### 3. The Proposed Two-Sample Empirical Likelihood Ratio Test for the Case-Control Study

In this section, we consider independent samples of sizes n and k from two populations. The data-points in each sample are independent and identically distributed. Let  $X_1, ..., X_n$  present a control sample from distribution  $F_X$  with a density function  $f_X(x)$ , and  $Y_1, ..., Y_k$  be a case sample from distribution  $F_Y$  with a density function  $f_Y(y)$ . We want to test the null hypothesis

$$H_0: F_Y = F_X = F_0 \text{ versus } H_1: F_Y \neq F_X = F_0,$$
 (7)

where distributions  $F_0 = F_X$  and  $F_Y$  are completely unknown. In the context of (7), the likelihood ratio test statistic is

$$\frac{\prod_{i=1}^{n} f_{X}(X_{i})\prod_{i=1}^{k} f_{Y}(Y_{i})}{\prod_{i=1}^{n} f_{X}(X_{i})\prod_{i=1}^{k} f_{X}(Y_{i})} = \prod_{i=1}^{k} h(Y_{i}) = \prod_{i=1}^{k} h(Y_{i}) = \prod_{i=1}^{k} h_{i}, \qquad (8)$$

where  $h_i = h(Y_{(i)}) = f_Y(Y_{(i)}) / f_X(Y_{(i)})$ , and  $Y_{(1)} \le Y_{(2)} \le ... \le Y_{(k)}$  are the order statistics based on the observations  $Y_1, ..., Y_k$ . (One can present  $f_Y = f_X h$ , where  $h = f_Y / f_X$ , and hence h can be considered as an unknown function under  $H_1$ .) Following the maximum EL methodology presented by Vexler and Gurevich [11], we find that values of  $h_i$ , i = 1, ..., k that maximize (8), satisfying some empirical constraints caused by the equation  $\int_{-\infty}^{+\infty} f_Y(u) du = \int_{-\infty}^{+\infty} f_X(u) h(u) du = 1$  are

 $h_{j} = 2m/(k(F_{Xn}(Y_{(j+m)}) - F_{Xn}(Y_{(j-m)}))), \quad j = 1,...,k \text{, where } F_{Xn}(x) = n^{-1} \sum_{i=1}^{n} I(X_{i} \le x) \text{ is the empirical distribution function } (I(\cdot) \text{ is the indicator function}). Here \quad Y_{(j)} = Y_{(1)}, \text{ if } j \le 1, \text{ and } Y_{(j)} = Y_{(k)},$ 

if  $j \ge k$ . Therefore, the maximum EL method yields the entropy-based test-statistic  $\prod_{j=1}^{k} \left[ \frac{2m}{k} \left( F_{Xn} \left( Y_{(j+m)} \right) - F_{Xn} \left( Y_{(j-m)} \right) \right) \right]$ . Finally, utilizing arguments of Section 1, we suggest the test-statistic

$$V_{nk} = \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k \left( F_{Xn} \left( Y_{(j+m)} \right) - F_{Xn} \left( Y_{(j-m)} \right) \right)}, \ 0 < \delta < 1$$
(9)

for the case-control problem (7).

The proposed test is to reject the null hypothesis of (7) if

 $\log(V_{nk}) > C , \qquad (10)$ 

where C is a test-threshold. (Similarly to Canner [3], we will arbitrarily define  $F_{Xn}(x) - F_{Xn}(y) = 1/(n+k)$ , if  $F_{Xn}(x) = F_{Xn}(y)$ .)

Significance level of the proposed test. Since  $I(X > Y) = I(F_0(X) > F_0(Y))$ , where  $F_0(x)$  is the cumulative distribution function of the distribution  $F_0$ , the significance level of the test (10) is  $P_{H_0} \{ \log(V_{nk}) > C \} = P_{X_1, \dots, X_n, Y_1, \dots, Y_k \sim Unif(0,1)} \{ \log(V_{nk}) > C \}$ . That is, the type I error of the proposed test (10) can be calculated exactly or approximately by simulations, for all sample sizes n, k and  $0 < \delta < 1$ . Fix  $\delta = 0.5$  in (9). Table 1 displays Monte Carlo roots C of the equations  $P_{X_1, \dots, X_n, Y_1, \dots, Y_k \sim Unif(0,1)} \{ \log(V_{nk}) > C \} = \alpha$ , for different values of  $\alpha$  and n, k. For each value of  $\alpha$ , n, k, the type I error results were derived via 55,000 generations of statistic  $\log(V_{nk})$ 's values.

**Table 1.** Critical values C for the test (10) with  $\delta = 0.5$ 

k												
	1	10	15	20	25	30	35	40	50	60	80	100
n												
	α											
10	0.01	9.704	10.482	11.797	13.213	14.751	16.459	18.255	22.167	26.409	35.735	45.594
	0.025	8.318	9.384	10.683	11.981	13.483	15.243	16.868	20.663	24.848	34.113	43.889
	0.05	7.507	8.468	9.560	11.000	12.384	13.857	15.770	19.500	23.462	32.439	42.503
	0.1	6.526	7.592	8.619	9.881	11.285	12.640	14.384	17.996	22.075	30.530	40.480
15	0.01	10 131	11 155	12 306	13/138	1/1 805	16 196	17 623	20.853	24 257	31.870	40 314
15	0.025	8 935	9 992	11 199	12 221	13 483	14 813	16 186	19 201	22 465	30.024	38 437
	0.025	8.060	9 181	10 190	11 240	12 503	13 609	15 024	17 915	21.038	28 466	36 773
	0.1	7.038	8.062	9.128	10.090	11.222	12.316	13.698	16.429	19.546	26.738	34.853
				,						-,		
20	0.01	10.397	11.694	12.934	14.083	15.261	16.2465	17.526	20.188	23.094	29.707	37.090
	0.025	9.246	10.456	11.666	12.731	13.939	14.899	16.086	18.683	21.420	27.898	35.069
	0.05	8.266	9.427	10.676	11.619	12.796	13.723	14.923	17.382	19.980	26.306	33.337
	0.1	7.148	8.228	9.468	10.402	11.519	12.430	13.565	15.841	18.328	24.390	31.309
25	0.01	10.589	12.039	13.271	14.438	15.699	16.545	17.834	20.205	22.610	28.471	34.975
	0.025	9.335	10.688	11.922	12.934	14.246	15.052	16.357	18.592	20.971	26.419	32.922
	0.05	8.254	9.489	10.773	11.764	13.014	13.814	15.085	17.202	19.481	24.802	30.998
	0.1	7.107	8.203	9.464	10.413	11.630	12.414	13.593	15.564	17.746	22.897	28.919
30	0.01	10.645	11.884	13.374	14.542	15.846	16.928	18.120	20.260	22.485	27.820	33.479
	0.025	9.380	10.466	11.881	13.001	14.284	15.274	16.447	18.603	20.719	25.810	31.332
	0.05	8.263	9.363	10.715	11.730	12.961	13.883	15.059	17.119	19.117	24.010	29.433
	0.1	7.083	8.083	9.375	10.265	11.458	12.326	13.443	15.427	17.356	21.960	27.133
35	0.01	10 594	11 915	13 306	14 345	15 826	16 731	18 091	20 194	22 413	27 187	32 425
55	0.025	9 196	10 373	11 789	12 748	14 067	15.096	16 331	18 4 59	20.483	25 140	30 123
	0.05	8.100	9.248	10 494	11 453	12.701	13.575	14.820	16.819	18 792	23.285	28.161
	0.1	6.953	7.909	9.154	9.980	11.178	11.973	13.149	15.045	16.848	21.132	25.794
						K						
----------	-------	--------	---------	--------	--------	--------	--------	--------	--------	--------	--------	--------
		10	15	20	25	30	35	40	50	60	80	100
n	~											
	u											
40	0.01	10.542	11.692	13.143	14.378	15.653	16.659	17.875	20.127	22.212	26.927	31.777
	0.025	9.140	10.174	11.649	12.633	13.933	14.749	16.116	18.112	20.171	24.653	29.379
	0.05	8.057	9.032	10.367	11.253	12.447	13.211	14.537	16.407	18.316	22.679	27.324
	0.1	6.880	7.782	8.994	9.798	10.882	11.606	12.742	14.490	16.272	20.312	24.705
50	0.01	10.250	11 5/18	12 860	13 744	15 137	15 875	17 3/3	19/66	21 533	26.053	30 3/9
50	0.025	8 924	9 997	11 241	12 013	13 282	14 045	15 403	17.400	19 242	20.055	27 738
	0.025	7 823	8 802	0.067	10 709	11 874	12 540	13.754	15 482	17 281	21 223	27.730
	0.05	6.678	7.607	8 655	9367	10 325	11.015	12 045	13 571	15 204	18 738	23.320
	0.1	0.078	7.007	0.055	2.307	10.525	11.015	12.045	15.571	15.204	10.750	22.400
60	0.01	10.312	11.157	12.501	13.428	14.524	15.238	16.559	18.439	20.254	24.498	28.806
	0.025	8.861	9.745	10.894	11.757	12.784	13.448	14.664	16.261	17.972	21.796	25.743
	0.05	7.737	8.660	9.669	10.411	11.371	12.019	13.112	14.562	16.207	19.684	23.256
	0.1	6.624	7.466	8.405	9.092	9.988	10.577	11.486	12.843	14.318	17.395	20.587
00	0.01	0.090	10.951	11 022	12 (42	12 020	14 075	15 202	17.002	10 (04	21.004	25 702
80	0.01	9.989	10.851	11.833	12.042	13.838	14.275	12.393	17.002	16.084	21.884	25.705
	0.025	8.620	9.457	0.269	11.120	12.134	12.085	13.0//	12.507	10.5/5	19.383	22.919
	0.05	/.5/0	8.330	9.208	9.893	10.815	11.300	12.231	13.397	12.220	17.708	20.792
	0.1	6.494	7.199	8.037	8.005	9.506	9.980	10.796	12.037	13.238	15.767	18.479
100	0.01	9.824	10.569	11.567	12.218	13.171	13.806	14.709	16.179	17.356	20.413	23.243
	0.025	8.499	9.221	10.187	10.724	11.621	12.262	13.017	14.324	15.537	18.149	21.008
	0.05	7.512	8.190	9.091	9.597	10.469	10.989	11.712	12.919	14.087	16.540	19.067
	0.1	6.448	7.085	7.892	8.446	9.205	9.692	10.349	11.486	12.549	14.866	17.043
$\infty$	0.01	9.138	9.623	10.291	10.650	11.113	11.449	11.888	12.507	13.125	14.232	15.225
	0.025	7.967	8.486	9.081	9.493	9.930	10.271	10.689	11.311	11.891	12.993	13.955
	0.05	7.029	7.579	8.157	8.555	9.017	9.343	9.722	10.363	10.937	12.004	12.933
	0.1	6.080	6.626	7.182	7.581	8.025	8.360	8.722	9.357	9.920	10.957	11.850

The following propositions present asymptotic operating characteristics of the test (10).

**Proposition 3.1.** For each  $0 < \delta < 1$ ,

$$V_{nk} \xrightarrow{P} \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k(Z_{(j+m)} - Z_{(j-m)})}, \text{ as } n \to \infty$$

where under  $H_0$ ,  $Z_1$ ,..., $Z_k \sim Unif(0,1)$ , whereas under  $H_1$ ,  $Z_1$ ,..., $Z_k \sim F_Z$ , and  $F_Z$  is a nonuniform distribution function with a density function  $f_Z$  concentrated on [0,1].

*Proof.* We note that, for each  $1 \le i \le k$  and  $\varepsilon > 0$ , we have

$$P(|F_{X_n}(Y_i) - F_X(Y_i)| > \varepsilon) = \int_{-\infty}^{\infty} P(|F_{X_n}(y) - F_X(y)| > \varepsilon) f_Y(y) dy \xrightarrow[n \to \infty]{} 0.$$

Therefore,  $F_{Xn}(Y_i) \xrightarrow{P} F_X(Y_i) = Z_i$ , as  $n \to \infty$ . Obviously, under  $H_0$ ,  $Z_i$  has the uniform Unif(0,1) distribution. Under  $H_1$ , the distribution of the  $Z_i$  is not uniform but concentrated on the interval [0,1].

When the control-sample size  $n \rightarrow \infty$ , the rule (10) rejects  $H_0$  if

$$\min_{1 \le m < k^{1-\delta}} \log \left( \prod_{j=1}^{k} \frac{2m}{k \left( Z_{(j+m)} - Z_{(j-m)} \right)} \right) > C$$

$$\tag{11}$$

(*C* is from (10)). Note that, (11) is an empirical likelihood modification of the well-known test for uniformity proposed by Dudewicz & Van Der Meulen [4]. Monte Carlo critical values *C* for the test (11), corresponding to different values of  $\alpha$  and *k*, are presented in the last lines of the Table 1 (signed  $n = \infty$ ). These critical values can be used for the two-sample test (10) based on data with a large number of controls  $X_1, ..., X_n$ .

**Proposition 3.2.** For each  $\delta \in (0,1)$ , under  $H_0$ ,  $k^{-1} \log(V_{nk}) \xrightarrow{P} 0$ , as  $n \to \infty$ ,  $k \to \infty$ ; under  $H_1$ ,  $k^{-1} \log(V_{nk}) \xrightarrow{P} b$ , as  $n \to \infty$ ,  $k \to \infty$ , where *b* is a positive constant.

*Proof.* By virtue of Proposition 3.1, for each  $0 < \delta < 1$ ,

$$k^{-1}\log(V_{nk}) \xrightarrow{P} k^{-1}\log\left(\min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k(Z_{(j+m)} - Z_{(j-m)})}\right), \text{ as } n \to \infty,$$

$$(12)$$

where  $Z_i = F_X(Y_i)$ . Let  $F_Z$  define the distribution of  $Z_j$  with a density function  $f_Z$ , j = 1,...,k. We pointed out in the proof of Proposition 3.1 that under  $H_0$ ,  $F_Z$  is the uniform Unif(0,1) distribution. Under  $H_1$ ,  $F_Z$  is not uniform but  $F_Z$  is concentrated on the interval [0,1]. Consider a behavior of

the statistic 
$$k^{-1} \log \left( \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} 2m \left( k \left( Z_{(j+m)} - Z_{(j-m)} \right) \right)^{-1} \right)$$
, as  $k \to \infty$ . Note tha  
 $k^{-1} \log \left( \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k \left( Z_{(j+m)} - Z_{(j-m)} \right)} \right) = -\max_{1 \le m < n^{1-\delta}} p_{mk}$ ,

where  $p_{mk} = k^{-1} \sum_{j=1}^{k} \log(k(Z_{(j+m)} - Z_{(j-m)})/2m)$ . Following Vasicek [10], after some reorganization, we obtain

$$p_{mk} = (2m)^{-1} \sum_{i=1}^{2m} S_i + U_{mk}, S_i = -\sum_{j=1}^k \log \left( \frac{F_Z(Z_{(j+m)}) - F_Z(Z_{(j-m)})}{Z_{(j+m)} - Z_{(j-m)}} \right) \left( F_{Zk}(Z_{(j+m)}) - F_{Zk}(Z_{(j-m)}) \right),$$
  
$$j \equiv i \pmod{2m}, U_{mk} = \frac{1}{k} \sum_{j=1}^k \log \left( \frac{k}{2m} \left( F_Z(Z_{(j+m)}) - F_Z(Z_{(j-m)}) \right) \right),$$

where  $F_Z(x)$  and  $F_{Zk}(x) = \frac{1}{k} \sum_{i=1}^k I(Z_i \le x)$  are the cumulative and empirical distribution functions, respectively. Vasicek [10] showed that  $S_i$  uniformly converges in probability to the entropy of the density  $f_Z(x)$  (as  $k \to \infty$ ,  $m/k \to 0$ ), for all  $1 \le m \le k^{1-\delta}$ ,  $0 < \delta < 1$ . The statistic  $U_{mk}$  is a non-positive random variable distributed independently of  $F_Z$  and  $U_{mk} \xrightarrow{P} 0$  as  $k \to \infty$ ,  $m \to \infty$ . Thus,

$$k^{-1} \log \left( \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k(Z_{(j+m)} - Z_{(j-m)})} \right) \le -p_{k^{1-\delta_k}} \xrightarrow{P} -H(f_Z), \text{ as } k \to \infty,$$
  
$$k^{-1} \log \left( \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k(Z_{(j+m)} - Z_{(j-m)})} \right) \ge -\max_{1 \le m < k^{1-\delta}} (2m)^{-1} \sum_{i=1}^{2m} S_i \xrightarrow{P} -H(f_Z), \text{ as } k \to \infty.$$

Therefore,

$$k^{-1} \log \left( \min_{1 \le m < k^{1-\delta}} \prod_{j=1}^{k} \frac{2m}{k \left( Z_{(j+m)} - Z_{(j-m)} \right)} \right) \xrightarrow{P} -H(f_Z), \text{ as } k \to \infty.$$
(13)

Since, for each density function f(x) concentrated on [0,1], one always has  $H(f) \le 0$ , with the maximum value H(f) = 0, being uniquely attained by the uniform Unif(0,1) density (Dudewicz & Van Der Meulen [4]), the equations (12) and (13) complete the proof.

Proposition 3.2 declares a consistency of the test (10) as  $n \to \infty$ ,  $k \to \infty$ .

### 4. Monte Carlo Study

In this section, we investigate the power properties of the proposed test (10) (with  $\delta = 0.5$ ) comparing with the commonly used two-sample Kolmogorov–Smirnov (KS) test (Birnbaum and Hall [2]; Massey [6]). To evaluate the properties of test (10), we conduct the following Monte Carlo simulations. For different values of n, k and different null/alternative distributions, 25,000 pairs of samples were generated corresponding to the problem (7). The test-statistics  $\log(V_{nk})$  with  $\delta = 0.5$  and the statistic of KS were calculated from each pair of the samples. The simulated powers of the tests are shown in Table 2, at the  $\alpha = 0.05$  level of significance. Table 2 does not display results of all simulations that we executed. We balance situations, in which the test (10) or the KS test can be recommended.

Baseline Distribution $F_X = F_0$	Alternative Distribution $F_{Y}$	n	k	Proposed test (10)	KS test
Norm(0,1)	Unif (-1,1)				
		45	45	0.9889	0.1490
		25	25	0.8304	0.0640
		15	25	0.7132	0.0878
		25	15	0.6059	0.0603
		15	15	0.5177	0.0363
Exp(1)	Lognorm(0,1)				
		45	45	0.4027	0.3362
		25	25	0.2074	0.1722
		15	25	0.1790	0.1696
		25	15	0.1406	0.1634
		15	15	0.1332	0.0866
		10	10	0.0964	0.1167
Norm(0,1)	Norm(0.5,1)				
		45	45	0.3755	0.5129
		15	15	0.1156	0.1399
Norm(0,1)	$Norm(0.5, 1.5^2)$				
		45	45	0.4098	0.5034
<i>Norm</i> (0,1)	$Norm(0, 1.5^2)$				
		50	50	0.2355	0.1248
		45	45	0.1669	0.1292
		25	25	0.0389	0.0663
		15	15	0.0142	0.0174
Exp(1)	<i>Exp</i> (1.5)				
		45	45	0.2670	0.3002
		25	25	0.1790	0.1582
		15	15	0.1257	0.0831

**Table 2.** The Monte Carlo powers of the proposed test (10) with  $\delta = 0.5$  and Kolmogorov–Smirnov (KS) test; a = 0.05.

Table 2 confirms that for relatively small and average sample sizes n and k the test (10) can be much more powerful than the KS test. (Table 2,  $F_X = Norm(0,1)$ ,  $F_Y = Unif(-1,1)$ ). In these cases with n = 15, k = 15, the power of the KS test is less than the type I error, whereas the Monte Carlo power of the proposed test is 0.5177. To explain this phenomenon, we would like to emphasize the known fact that entropy-based tests for goodness-of-fit are very powerful for detecting a change towards small variance (Dudewicz & Van Der Meulen [4]). It seems that the test (10) has also a high level of the power when the variance of the alternative distribution is smaller than that of the baseline distribution. The KS test sometimes demonstrates the powers that are better than those of the proposed test (10), for relatively small n and k (for example, when  $F_X = Norm(0,1)$ ,  $F_Y = Norm(0.5,1)$ ). However, for all sample sizes n and k, if the variance of the alternative distribution  $F_Y$  is smaller than that of the baseline distribution  $F_X$  then the entropy-based test (10) is superior to the KS test.

# Conclusions

In this article, we have presented a methodology for developing density-based EL tests. The objective was to indicate there is a unified method to derive test-statistics based on samples entropy. This method utilizes the main idea of the empirical likelihood methodology where the empirical

likelihood function consists of components that maximize this likelihood function and satisfy empirical constraints. We have proved that the entropy-based tests for goodness-of-fit have the EL ratio structure. We focused on the tests for normality and uniformity. The test for exponentiality based on Kullback–Leibler information (Ebrahami et. al [5]) can be also considered paying attention to Section 1. In addition, we applied the proposed density-based EL ratio technique to create two-sample EL test for the case-control study based on samples entropy. The presented Monte Carlo simulations confirm that the proposed nonparametric test is superior to the standard procedures. While considering the approach of this article, nonparametric two-sample entropy based tests for a general case as well as k-sample entropy-based tests can be easily constructed. We believe that the proposed approach can be applied to create a nonparametric two-sample test for detecting shift alternatives. (In this case, the two-sample Mann–Whitney–Wilcoxon test is a common procedure.) Further studies are needed to test the suggested approach in other contexts. We hope that this article will stimulate future theoretical and applied research on this topic.

# Acknowledgements

This work was partially supported by the Internal Funding Program of the Shamoon College of Engineering (SCE).

# References

- 1. Arizono, I., Ohta, H. A test for normality based on Kullback–Leibler information, *The American Statistician*, Vol. 43, 1989, pp. 20–22.
- Birnbaum, Z. W., Hall, R. A. Small sample distributions for multi-sample statistics of the Smirnov type, Annals of Mathematical Statistics, Vol. 31, 1960, pp. 710–720.
- 3. Canner, P. L. A simulation study of one-and two-sample Kolmogorov–Smirnov statistics with a particular weight function, *Journal of the American Statistical Association*, Vol. 70, 1975, pp. 209–211.
- 4. Dudewicz, E. J., van der Meulen, E. C. Entropy-based tests of uniformity, *Journal of the American Statistical Association*, Vol. 76, 1981, pp. 967–974.
- 5. Ebrahami, N., Habibullah, M., Soofi, E. S. Testing Exponentiality Based on Kullback–Leibler Information, *Journal of the Royal Statistical Society. Series B.*, Vol. 54, 1992, pp. 739–748.
- 6. Massey, F. The distribution of the maximum deviation between two sample cumulative step functions, *Annals of Mathematical Statistics*, Vol. 22, 1951, pp. 125–128.
- 7. Owen, A. B. *Empirical likelihood*. Chapman & Hall/CRC, 2001.
- 8. Park, S., Park, D. Correcting moments for goodness of fit tests based on two entropy estimates, *Journal of Statistical Computation and Simulation*, Vol. 73, 2003, pp. 685–694.
- 9. Tusnady, G. On asymptotically optimal tests. Annals of Statistics, Vol. 5, 1977, pp. 385–393.
- 10. Vasicek, O. A test for normality based on sample entropy, *Journal of the Royal Statistical Society*, *Series B.*, Vol. 38, 1976, pp. 54–59.
- 11. Vexler, A., Gurevich, G. Empirical likelihood ratios applied to goodness-of-fit tests based on sample entropy, *Computational Statistics and Data Analysis*, Vol. 54, pp. 531–545.
- 12. Zhang, J. Powerful goodness-of-fit tests based on the likelihood ratio, *Journal of the Royal Statistical Society. Series B.*, Vol. 64, 2002, pp. 281–294.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 40–49 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# ON THE ESTIMATION OF STRUCTURAL PARAMETERS IN FRAILTY MODELS FOR INTERVAL CENSORED AND TRUNCATED DATA

# **F.** Vonta<sup>1</sup>, C. Huber<sup>2</sup>

 <sup>1</sup> Department of Mathematics National Technical University of Athens, Athens, Greece E-mail: vonta@math.ntua.gr
 <sup>2</sup> Universite Rene Descartes-Paris 5 45 rue des Saints-Peres, 75006 Paris, France E-mail: huber@paris.descartes.fr

We consider survival data that are both interval censored and interval truncated. We assume a semiparametric frailty or transformation model for the survival function and consider censoring and truncation distributions as in Huber, Solev and Vonta [6], [7]. We propose the use of modified profile likelihood estimators for the structural parameter of the model as in Slud and Vonta [11]. For fixed values of the structural parameter, we derive the least favourable parametrization of the nuisance infinite-dimensional parameter, on which the definition of the modified profile likelihood estimator is relied upon. We discuss the semiparametric efficiency of the modified profile likelihood estimator of the finite-dimensional nuisance parameter, that is, the baseline cumulative hazard function.

Keywords: frailty models, least favourable model, interval censored and truncated data, semiparametric estimation

#### **1. Introduction**

Many times we are faced with complex observational schemes such as interval censored and interval truncated data. For example, HIV infection or toxicity of a treatment, is not exactly known, but it is usually known to have taken place between two dates  $t_1$  and  $t_2$ . Furthermore, some people may be lost

from the sample if they are observed during a period of time not including such pair of dates  $t_1, t_2$ .

Turnbull [12] proposed a method for nonparametric maximum likelihood estimation of the distribution function in the case of arbitrarily censored and truncated data. His method, slightly corrected by Frydman [5], has been used extensively since by several authors, and extended to the Cox model by Alioum and Commenges [1] and to the frailty or transformation models by Huber and Vonta [8]. In Huber, Solev and Vonta [6] and [7] we give conditions on the involved distributions, namely, the censoring, truncation and survival distributions, implying the consistency of a nonparametric maximum likelihood estimator of the density of the survival process in the nonparametric case. We also provide the rate of convergence of the NPMLE of the density within a certain class of density functions.

In Bickel et al. [2] and van der Vaart and Wellner [13] one can find several different tools for expressing the semiparametric information about the finite-dimensional parameter of interest in semiparametric models and for establishing semiparametric efficiency. Slud and Vonta [11] proposed modified profile likelihood estimators for the finite-dimensional parameter of interest in the presence of the infinite-dimensional nuisance parameter. These estimators have been shown to be semiparametric efficient under regularity conditions. Results in Slud and Vonta [11] generalize those of Severini and Wong [10] in the semiparametric case.

In this paper we propose to apply the modified profile likelihood approach to the case of censored and truncated data assuming a transformation semiparametric model (Vonta [15]) for the survival time. The transformation models include the frailty models which arise as a generalization of the Cox model (Cox [4]) when one introduces a random effect term into it, with the purpose of explaining possible population heterogeneity which remains unexplained from the Cox model. The class of transformation models that we consider are equivalent to the class of models defined in Cheng, Wei and Ying [3].

In section two, we give a representation of the censoring and truncation mechanisms. In section three, we define the proposed frailty or transformation model as well as the modified profile likelihood estimators. We also establish the form of the least favourable nuisance parametrization on which the modified profile likelihood approach is relied upon. The least favourable parametric submodel cannot be given in closed form. We derive a recursive equation through which the least favourable model is implicitly defined.

Finally, we discuss the semiparametric efficiency of the modified profile likelihood estimator of the regression parameter  $\beta$ .

# 2. The Observation Scheme

Time X to an event that changes permanently the state of subject i under study (state 0 before X, 1 afterwards) is a random variable whose distribution is to be estimated under the following observation scheme:

1. Censorship: observation of each subject i does not take place continuously but is scheduled at a (random) number K(i) of (random) times

$$a < Y_{i,1} < \cdots < Y_{i,K(i)} < b$$

where usually a will be equal to 0 and b is a finite strictly positive number. Let  $\tau_i := \{Y_{i,j}, j = 1, \dots, K(i)\}$  be the set of scheduled observation times for subject i and  $t_i := \{y_{i,j}, j = 1, \dots, K(i)\}$  a realization of  $\tau_i$ .

2. Truncation: only those elements of  $t_i$  that are inside a given (random) truncating window  $(Z_{i,1} Z_{i,2}]$  give rise to an actual observation of subject i.

Thus, if subject i is observed in state 0 at time  $y_{i,j}$  and in state 1 at time  $y_{i,j+1}$ , inside its window  $\Delta := (z_{i,1}, z_{i,2}]$ , one observes subject i at all times of  $t_i$  included in  $\triangle$ . A sufficient statistic for this problem is thus the two embedded intervals "bracketing" the unobserved X = x:

 $z_{i,1} \le y_{i,k_1} \le y_{i,j} < x \le y_{i,j+1} \le y_{i,k_2} \le z_{i,2}$ ,

where  $y_{i,k_1}$  is the smallest time in  $t_i$  which is greater than or equal to  $z_{i,1}$  and  $y_{i,k_2}$  is the largest time in  $t_i$  that is less than or equal to  $z_{i,2}$ .

# 2.1. Censoring

Let  $\tau$  be a random partition defined on ]a,b], where usually a will be equal to 0 and b is a finite strictly positive number:

$$\tau = \{Y_0 = a < Y_1 < \ldots < Y_K < Y_{K+1} = b, \bigcup_{j=0}^K (Y_j, Y_{j+1}] = (a, b]\},$$
(1)

where K is a fixed number or a random number with known law in  $\{2, ..., K_0\}$  for some given  $K_0$  such that  $2 < K_0 < \infty$ .

For each  $x \in (a, b)$  we define

$$j(x) = \inf\left\{j : x \le Y_{j+1}\right\},\tag{2}$$

$$\mathcal{G}(x) = \left(Y_{\mathbf{j}(x)}, Y_{\mathbf{j}(x)+1}\right] \coloneqq \left(L(x), R(x)\right], \ x \in (a,b),$$
(3)

where L(x) and R(x) may be thought of as the left and right values in partition  $\tau$  that "bracket" (the survival) X = x.

Then it is clear that

$$\vartheta(x) = \vartheta(y), \text{ or } \partial(x) \cap \partial(y) = \emptyset$$
(4)

and we call  $\mathcal{G}(x)$  a simple random covering of (a,b). From now on we will take a to be 0 for convenience.

## 2.2. Truncation

Let  $\mathcal{G}(x) = (L(x), R(x)], x \in \mathbb{R}$ , be the simple random covering defined by the partition  $\tau = t := \{y_j, j = 1, 2, ..., k\}$ . Then, a fixed interval  $\Delta = (z_1, z_2]$ , and z the associated vector  $(z_1, z_2), z_1 \le z_2$ , is a truncating interval. This means that the only available observations of the subject i under investigation take place at times that are those elements of t that are included in  $(R(z_1); L(z_2)]$ , which behaves like the effective "truncating window".

Our basic assumption is the following: First we assume that the random covering  $\mathcal{G}(\cdot)$ , the random variable X and the random interval  $\Delta$  are independent. Second, we assume that the distribution of any  $(Y_{m+1}, \dots, Y_n)$  is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^{n-m}$  and that the distribution of  $(Z_1, Z_2)$  is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^2$ .

In that case a summary of the observations on the subject under investigation is the pair of embedded intervals:

$$R(z_1) \le L(x) < R(x) \le L(z_2),$$

where the censoring interval (L(X), R(X)] of the covering  $\mathcal{G}(\cdot)$ , contains X, and the random interval  $\Delta^* = (R(z_1), L(z_2)]$  actually truncates X.

When  $(L(X), R(X)] \not\subset \Delta$  we do not have any observation.

In the special case of right truncation

 $\Delta = (0, z]$ 

and a summary of the observations on the subject under investigation is the triple of random variables (L(x), R(x), L(z)) such that:

 $0 \le L(x) < R(x) \le L(z).$ 

#### 2.3. Combined Censoring and Truncation

Let us define

1) Conditionally on a fixed value t of  $\tau$  the random interval  $\Delta$  is taken from the truncating distribution

 $\mathcal{P}_t \{A\} = P\{\Delta \in A | \text{ the interval } (\mathbb{Z}_1, \mathbb{Z}_2] \text{ contains at least two points of } t\}.$ 

In other words, conditionally on fixed values of  $\tau = t$  the random vector  $Z = (Z_1, Z_2)$  is taken from the truncating distribution

$$P_{t}\{B\} = P\{Z \in B | R(Z_{1}) < L(Z_{2})\}$$

2) Conditionally on a fixed value of  $\tau = t$  and  $\Delta = \Delta = (z_1, z_2]$ , the random variable X is taken from the truncating distribution

$$P_{t, \perp} \{C\} = P\{X \in C | X \in (R(z_1), L(z_2)]\}.$$

In other words conditionally on fixed values of  $\tau = t$  and  $Z_1 = z_1, Z_2 = z_2$  the random variable X is taken from the truncating distribution

$$P\{C|t, z_1, z_2\} = P\{X \in C | X \in (R(z_1 | t), L(z_2 | t))\}.$$
(5)

We consider now the simple case of right truncation by Z, where for a random variable Z the random interval  $\Delta = (0, Z]$ . We denote for short when there is no ambiguity about the partition  $\tau = t$  simply:

$$L(Z) \coloneqq L(Z \mid \tau = t).$$

Then, conditionally on fixed values of  $\tau = t$  and Z = z the random variable X is taken from the truncating distribution

$$P\left\{C \mid t, z\right\} = P\left\{X \in C \mid X \leq L(z)\right\}.$$

# 3. The Model

## 3.1. Definition

Let X be a random variable with density f and survival function conditional on  $\Xi = \xi$  defined by (Vonta [15])

$$S(t \mid \Xi = \xi) = P(X > t \mid \Xi = \xi) = e^{-G(e^{\beta^{t}} \xi_{\Lambda(t)})},$$
(6)

where  $\beta \in \mathbb{R}^d$  is the parameter of interest,  $\Xi$  a d-dimensional vector of covariates,  $\Lambda$  the baseline cumulative hazard function which plays the role of the infinite-dimensional nuisance parameter and G a known function which satisfies some regularity conditions such as G increasing and concave with G(0) = 0 and  $G(\infty) = \infty$ . The notation  $\beta^T$  denotes from now on, the transpose of the vector  $\beta$ .

Let  $v^k$  be the Lebesgue measure on  $\mathbb{R}^k$  and recall that  $\xi \in \mathbb{R}^d$  and let  $dP_{\Xi} = \phi(\xi) dv_d$  for some  $\sigma$  – finite measure  $v_d$  on  $\mathbb{R}^d$  possibly equal to  $v^d$ . Without loss of generality we consider the right-truncation case. The problem that we are faced with could be formulated as follows. Our observations are  $Q_1, \dots, Q_n$ , i.i.d. random vectors,  $Q = (L(X), R(X), L(Z), \Xi)$ , with density  $p(u, v, w, \xi)$  with respect to a measure  $v^*$  (Huber, Solev and Vonta [6]) given as

$$p(u, v, w, \xi) = r(u, v, w) \cdot \frac{\int_{0}^{v} f(t \mid \xi) dt}{\int_{0}^{w} f(t \mid \xi) dt} \cdot \phi(\xi),$$

$$(7)$$

where  $v^*$  is the measure on  $\mathbb{R}^{3+d}$  which is defined for continuous nonnegative functions  $\psi(s) = \psi(u, v, w, \xi)$  by the relation

$$\iiint \psi(s) dv^* = \iiint \psi(u, v, w, \xi) d(v^3 \otimes v_d)(u, v, w, \xi) + \iiint \psi(u, v, v, \xi) d(v^2 \otimes v_d)(u, v, \xi).$$

From model (6) p is equal to

$$r(u, v, w) \cdot \frac{\int_{0}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt}{\int_{0}^{w} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt} \phi(\xi)$$

$$\equiv r(u, v, w) \cdot \varphi(u, v, w, \xi \mid \beta, \lambda) \cdot \phi(\xi)$$
(8)

thus defining function  $\varphi$ . Here  $0 \le u < v \le w \le b$ , *r* is the known density with respect to  $v^{**}$ , the marginal of  $v^*$  integrated over  $\xi$ , of the known joint law of censoring and truncation. Function  $\varphi$  is the likelihood of each unobserved survival conditional on the censoring, truncation and covariate.

Recall that  $\beta$  is the parameter of interest and  $\Lambda$  the baseline cumulative hazard function, or equivalently  $\lambda$  the baseline hazard intensity function, is the infinite-dimensional nuisance parameter.

The joint law of censoring and truncation with density r with respect to  $v^{**}$  has two components, one denoted by  $r_3$  which is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^3$  (corresponding to the case where u < v < w) and a second one, denoted by  $r_2$  which is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^2$  (corresponding to the case where u < v = w). For details and an example of such a law r see Huber, Solev and Vonta [6].

We are interested in the efficient estimation of the parameter of interest  $\beta$  in the presence of the unknown cumulative hazard function  $\Lambda$  or equivalently in the presence of the hazard intensity function  $\lambda$  where obviously  $\lambda(t) \ge 0$ . We assume that the function  $\Lambda(t)$  is finite for finite times t and that  $\Lambda(\infty) = \infty$ .

#### **3.2. Modified Profile Likelihood Estimators**

In the notation of Slud and Vonta [11], assume that the independent identically distributed (*iid*) data-sample  $X_1, X_2, ..., X_n$  of random vectors in  $\mathbb{R}^k$  is observed and assumed to follow a marginal probability law  $\mu = P_{(\beta_0, \lambda_0)}$  where  $\beta_0 \in U \subset \mathbb{R}^d$ ,  $\lambda_0 \in V \subset L^0(\mathbb{R}^q, \nu)$  (Borel-measurable functions), where U is a fixed open set; V is a fixed set of positive measurable functions; and the  $\sigma$ -finite measure  $\nu$  (locally finite, but not necessarily a probability measure) is fixed on  $\mathbb{R}^q$ .

In addition assume that there is a family  $\{P_{(\beta,\lambda)}, (\beta,\lambda) \in U \times V\}$  of Borel probability measures on  $\mathbb{R}^k$ , such that for all  $(\beta,\lambda) \in U \times V$ ,  $P_{(\beta,\lambda)} \ll \mu$ , and the regularity of densities  $f_X(\cdot,\beta,\lambda) \equiv dP_{(\beta,\lambda)}/d\mu$  as functions of  $(\beta,\lambda)$  are further restricted by assumptions given in detail in Slud and Vonta [11]. Note that by definition,  $f_X(\cdot,\beta_0,\lambda_0) \equiv 1$ . The true parameter-component value  $\beta_0$  is assumed to lie in the interior of a fixed, known compact set  $F \subset U$ .

The log-likelihood for the models  $P_{(\beta,\lambda)}$  and data  $\mathbf{X} = \{X_i\}_{i=1}^n$  is defined by  $\log lik_n(\beta,\lambda) = \sum_{i=1}^n \log f_X(X_i,\beta,\lambda), \quad (\beta,\lambda) \in (U \times V).$ 

Define the Kullback-Leibler functional by

$$K(\beta,\lambda) \equiv -\int \log f_X(x,\beta,\lambda) d\mu(x)$$

The key idea of *modified profile likelihood* (Severini and Wong [10]) is to replace the nuisance parameter  $\lambda$  in the log-likelihood by a suitable estimator  $\tilde{\lambda}_{\beta}$  of the minimizer  $\lambda_{\beta}$  over  $\lambda \in V$ , of the Kullback–Leibler functional. The minimizer is assumed to be unique, smooth in  $\beta$  and the estimator of the minimizer is also assumed to be smooth in  $\beta$  and consistent for  $\lambda_{\beta}$  (Slud and Vonta [11]). The modified profile likelihood, is then the maximizer with respect to  $\beta$ , of the modified profile likelihood,

$$\sum_{i=1}^{n} \log f_X(X_i, \beta, \tilde{\lambda}_{\beta})$$

proved in Slud and Vonta [11] to be semiparametric efficient under the assumed therein regularity conditions. The d-dimensional parametric submodel ( $\beta$ ,  $\lambda_{\beta}$ ) is called a least favourable parametric submodel for the general semiparametric model  $P_{\beta,\lambda}$ , where  $\lambda_{\beta}$  is the minimizer of the Kullback–Leibler functional with respect to  $\lambda$  for fixed  $\beta$ .

Keeping  $\beta$  fixed we differentiate with respect to the parameter  $\lambda$  in the sense of Ga teaux differentiation. We consider perturbations of functions  $\lambda \in V$  by small multiples of functions  $\gamma$  in subsets of  $G \subseteq G_0$ . All spaces V, G and  $G_0$  are described in Slud and Vonta [11] but they ought to be specified appropriately according to the situation.

Here and in what follows, we define the differentiation operator  $D_{\lambda}$  for all functionals  $\Phi: V \to \mathbb{R}$ , and all  $\gamma \in G_0$ , by:

$$(D_{\lambda} \Phi(\lambda))(\gamma) \coloneqq \frac{d}{d\theta} \Phi(\lambda + \theta \gamma) \Big|_{\theta=0} \coloneqq \frac{d}{d\theta} \Phi(\lambda_{\theta}) \Big|_{\theta=0} , \qquad (9)$$

where  $\theta$  belongs in a small neighborhood of 0 and  $\lambda_{\theta}$  is defined as

$$\lambda_{\theta}(t) = \lambda(t) + \theta \gamma(t)$$
 leading to  $\Lambda_{\theta}(t) = \Lambda(t) + \theta \int_{0}^{t} \gamma(s) ds = \Lambda(t) + \theta \Gamma(t).$ 

# **3.3. Least Favorable Model**

The data-space  $\mathbf{D} = \mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}^d$  consists of vectors  $s = (u, v, w, \xi)$ . We denote the true parameters by  $(\beta_0, \lambda_0)$ .

We define the probability law for the true model by

$$d\mu(s) \equiv p_0(s)dv^*(s),$$

where  $p_0(s)$  denotes the density p taken at the true point  $(\beta_0, \lambda_0)$ .

Taking the densities  $f_Q(s | \beta, \lambda)$  with respect to  $\mu$ , the true law of the observations Q, as in Slud and Vonta [11], we get that  $f_O(s | \beta, \lambda)$  is equal to

$$\frac{\int_{u}^{v} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi}\lambda(t)dt}{\int_{0}^{w} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi}\lambda(t)dt} \cdot \frac{\int_{0}^{w} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda_{0}(t)} e^{\beta^{T}\xi}\lambda_{0}(t)dt}{\int_{u}^{v} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda_{0}(t)} e^{\beta^{T}\xi}\lambda_{0}(t)dt}$$
(10)

Following the methodology of Slud and Vonta [11] we will find in this section, for fixed  $\beta$ , the least favourable parametric submodel ( $\Lambda_{\beta}, \beta$ ) of the proposed semiparametric model.

The Kullback-Leibler functional is given by

$$K(\beta,\lambda) = -\int \log \left( \frac{p(u,v,w,\xi \mid \beta,\lambda)}{p(u,v,w,\xi \mid \beta_0,\lambda_0)} \right) p(u,v,w,\xi \mid \beta_0,\lambda_0) dv^*.$$

Due to the form of the law r(u, v, w) the Kullback–Leibler functional is written equivalently as

$$-\int \log \left( \frac{p(u,v,w,\xi \mid \beta,\lambda)}{p(u,v,w,\xi \mid \beta_0,\lambda_0)} \right) p(u,v,w,\xi \mid \beta_0,\lambda_0) d(v^3 \otimes v_d)$$
  
$$-\int \log \left( \frac{p(u,v,v,\xi \mid \beta,\lambda)}{p(u,v,v,\xi \mid \beta_0,\lambda_0)} \right) p(u,v,v,\xi \mid \beta_0,\lambda_0) d(v^2 \otimes v_d),$$

which because of model (6) is equal to

$$-\int \left\{ \log \left( \int_{u}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt \right) -\log \left( \int_{0}^{w} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt \right) \right\} p(u,v,w,\xi \mid \beta_{0},\lambda_{0}) d(v^{3} \otimes v_{d}) \\ -\int \left\{ \log \left( \int_{u}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt \right) -\log \left( \int_{0}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))} G'(e^{\beta^{T}\xi}\Lambda(t)) e^{\beta^{T}\xi}\lambda(t) dt \right) \right\} p(u,v,v,\xi \mid \beta_{0},\lambda_{0}) d(v^{2} \otimes v_{d}) + C',$$
(11)

where  $C^{'}$  denotes a term that does not depend on  $(eta,\lambda)$  .

Let us define

$$p_{0,3} = p(u, v, w, \xi | \beta_0, \lambda_0)$$
  

$$p_{0,2} = p(u, v, v, \xi | \beta_0, \lambda_0).$$

The G  $\hat{\mathbf{a}}$  teaux differentiation of  $K(\beta, \lambda)$  in the direction  $\gamma$  is given as

$$\frac{d}{d\theta}K(\beta,\lambda_{\theta})|_{\theta=0} =
-\int \left\{ \frac{\int_{u}^{v} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi} \left( -G' + \frac{G''}{G'}|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi}\lambda(t)\Gamma(t) + \gamma(t) \right) dt}{\int_{u}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt} - \frac{\int_{0}^{w} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi} \left( -G' + \frac{G''}{G'}|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi}\lambda(t)\Gamma(t) + \gamma(t) \right) dt}{\int_{0}^{w} e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt} \right\} p_{0,3}d(v^{3}\otimes v_{d}) \\
-\int \left\{ \frac{\int_{u}^{v} e^{-G}G'|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi} \left( -G' + \frac{G''}{G'}|_{e^{\beta^{T}\xi}\Lambda(t)} e^{\beta^{T}\xi}\lambda(t)dt - \int_{0}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt - \int_{u}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt - \int_{u}^{v} e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt \right\} p_{0,2}d(v^{2}\otimes v_{d}).$$

$$(12)$$

By an integration by parts, the integral

$$\int_{u}^{v} \Gamma(t) d\left(e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\right) =$$
  
$$\Gamma(t)e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\Big|_{u}^{v} - \int_{u}^{v}e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\gamma(t)dt$$

and therefore the first numerator of (12) simplifies to

$$\Gamma(t)e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\Big|_{u}^{v}.$$

By a similar integration by parts we simplify the other three numerators of (12) to get that

$$\frac{d}{d\theta}K(\beta,\lambda_{\theta})|_{\theta=0} \text{ is equal to} 
-\int \left\{\frac{\Gamma(v)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(v)}e^{\beta^{T}\xi}-\Gamma(u)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(u)}e^{\beta^{T}\xi}}{\int_{u}^{v}e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt} 
-\frac{\Gamma(w)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(w)}e^{\beta^{T}\xi}-\Gamma(0)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(0)}e^{\beta^{T}\xi}}{\int_{0}^{w}e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt}\right\}p_{0,3}d(v^{3}\otimes v_{d}) 
-\int \left\{\frac{\Gamma(v)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(v)}e^{\beta^{T}\xi}-\Gamma(u)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(u)}e^{\beta^{T}\xi}}{\int_{u}^{v}e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt} 
-\frac{\Gamma(v)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(v)}e^{\beta^{T}\xi}-\Gamma(0)(e^{-G}G')|_{e^{\beta^{T}\xi}\Lambda(v)}e^{\beta^{T}\xi}}{\int_{0}^{v}e^{-G(e^{\beta^{T}\xi}\Lambda(t))}G'(e^{\beta^{T}\xi}\Lambda(t))e^{\beta^{T}\xi}\lambda(t)dt}\right\}p_{0,2}d(v^{2}\otimes v_{d}).$$

Finally, we have

$$\begin{aligned} &\frac{d}{d\theta} K(\beta, \lambda_{\theta})|_{\theta=0} = \\ &\int \Big\{ \frac{(\Gamma SG'(u) - \Gamma SG'(v))e^{\beta^{T}\xi}}{S(u) - S(v)} + \frac{\Gamma SG'(w)e^{\beta^{T}\xi}}{1 - S(w)} \Big\} p_{0,3} d(v^{3} \otimes v_{d}) \\ &+ \int \Big\{ \frac{(\Gamma SG'(u) - \Gamma SG'(v))e^{\beta^{T}\xi}}{S(u) - S(v)} + \frac{\Gamma SG'(v)e^{\beta^{T}\xi}}{1 - S(v)} \Big\} p_{0,2} d(v^{2} \otimes v_{d}) \end{aligned}$$

The survival function  $S(.) = S(.|\xi)$  and  $G'(.) = G'(.|\xi)$  but the dependence on  $\xi$  is omitted for convenience of notation.

Then we set the above derivative equal to 0 to obtain the equation

$$\int \left\{ \frac{\left( \Gamma SG'(u) - \Gamma SG'(v) \right) e^{\beta^{T}\xi}}{S(u) - S(v)} + \frac{\Gamma SG'(w) e^{\beta^{T}\xi}}{1 - S(w)} \right\} p_{0,3} d(v^{3} \otimes v_{d}) \\
+ \int \left\{ \frac{\left( \Gamma SG'(u) - \Gamma SG'(v) \right) e^{\beta^{T}\xi}}{S(u) - S(v)} + \frac{\Gamma SG'(v) e^{\beta^{T}\xi}}{1 - S(v)} \right\} p_{0,2} d(v^{2} \otimes v_{d}) = 0,$$
(13)

which should hold for all  $\gamma \in G_0$ . Equation (13) defines implicitly the minimizer  $\Lambda_\beta$  through which  $\lambda_\beta$  is subsequently defined.

**Lemma 1.** For the observational scheme defined in section section 2 under which  $Q = (L(X), R(X), L(Z), \Xi)$  is the observed random vector of variables with values  $s = (u, v, w, \xi)$ ,  $\mu$  the true law such that  $d\mu = p_0 dv^*$  and under model (6) assumed for the survival time X, for fixed  $\beta$ , a least favourable nuisance parametrization  $\Lambda_{\beta}$  is defined recursively through the equation

$$\Lambda_{\beta}(t) = \frac{I_1(t,b,\beta)}{I_2(t,b,\beta)},\tag{14}$$

where

$$\begin{split} &I_{1}(t,b,\beta) = -E_{\mu} \Biggl( \frac{\Lambda_{\beta}(L(X))S_{\beta}(L(X))G_{\beta}^{'}(L(X))e^{\beta^{T}\xi}}{S_{\beta}(L(X)) - S_{\beta}(R(X))} \,\Big| \, L(X) < t \Biggr) \\ &+ E_{\mu} \Biggl( \frac{\Lambda_{\beta}(R(X))S_{\beta}(R(X))G_{\beta}^{'}(R(X))e^{\beta^{T}\xi}}{S_{\beta}(L(X)) - S_{\beta}(R(X))} \,\Big| \, R(X) < t \Biggr) \\ &- E_{\mu} \Biggl( \frac{\Lambda_{\beta}(L(Z))S_{\beta}(L(Z))G_{\beta}^{'}(L(Z))e^{\beta^{T}\xi}}{1 - S_{\beta}(L(Z))} \,\Big| \, L(Z) < t \Biggr) \end{split}$$

and

$$\begin{split} &I_{2}(t,b,\beta) = E_{\mu} \Biggl( \frac{\Lambda_{\beta}(L(X))S_{\beta}(L(X))G_{\beta}^{'}(L(X))e^{\beta^{T}\xi}}{S_{\beta}(L(X)) - S_{\beta}(R(X))} \Big| L(X) \ge t \Biggr) \\ &- E_{\mu} \Biggl( \frac{\Lambda_{\beta}(R(X))S_{\beta}(R(X))G_{\beta}^{'}(R(X))e^{\beta^{T}\xi}}{S_{\beta}(L(X)) - S_{\beta}(R(X))} \Big| R(X) \ge t \Biggr) \\ &+ E_{\mu} \Biggl( \frac{\Lambda_{\beta}(L(Z))S_{\beta}(L(Z))G_{\beta}^{'}(L(Z))e^{\beta^{T}\xi}}{1 - S_{\beta}(L(Z))} \Big| L(Z) \ge t \Biggr). \end{split}$$

The least favourable direction  $\gamma$  is defined implicitly through equation (14) as  $\nabla_{\beta} \frac{d\Lambda_{\beta}(t)}{dt}|_{\beta=\beta_0}$ 

 $= \nabla_{\beta} \lambda_{\beta}(t) |_{\beta = \beta_0}.$ 

In order to establish semiparametric efficiency of the modified profile likelihood estimator of the parameter  $\beta$  we need to verify that all assumptions  $\mathbf{A}_0 - \mathbf{A}_8$  given in Slud and Vonta [11] are satisfied in the current situation. For this we will need to impose further regularity conditions on the function G and specify the spaces V, G and  $G_0$ . This task however is deferred to another paper.

# References

- 1. Alioum, A., Commenges, D. A proportional hazards model for arbitrarily censored and truncated data, *Biometrics*, Vol. 52, 1996, pp. 512–524.
- Bickel, P., Klaassen, C., Ritov, Y., Wellner, J. Efficient and Adaptive Inference in Semiparametric Models. Baltimore: Johns Hopkins Univ. Press, 1993.
- Cheng, S. C., Wei, L. J., Ying Z. Analysis of transformation models with censored data, *Biometrika*, Vol. 82, 1995, pp. 835–845.
- Cox, D. R. Regression models and life tables (with discussion), *Jour. Roy. Statist. Soc. Ser.* B, Vol. 34, 1972, pp. 187–202.
- 5. Frydman, H. A note on nonparametric estimation of the distribution function from interval-censored and truncated observations, *Journal of the Royal Statistical Society, Series B*, Vol. 56, 1994, pp. 71–74.
- Huber-Carol, C., Solev, V., Vonta, F. Interval censored and truncated data: rate of convergence of NPMLE of the density, *Journal of Statistical Planning and Inference*, Vol. 139, 2009, pp. 1734–1749.
- Huber-Carol, C., Solev, V., Vonta, F. Estimation of density for arbitrarily censored and truncated data. In: Probability, Statistics and Modelling in Public Health / M. S. Nikulin, D. Commenges, and C. Huber-Carol (Eds.). New York: Springer (Kluwer Acad. Publ.), 2006, pp. 246–265.

- 8. Huber-Carol, C., Vonta, F. Frailty models for arbitrarily censored and truncated data, *Lifetime Data Analysis*, Vol. 10, 2004, pp. 369–388.
- 9. Kosorok, M., Lee, B., Fine, J. Robust inference for univariate proportional hazards frailty regression models, *Ann. Statist.*, Vol. 32, No 4, 2004, pp. 1448–1491.
- 10. Severini, T. and W. Wong. Profile likelihood and conditionally parametric models, *Ann. Statist.*, Vol. 20, 1992, pp. 1768–1802.
- 11. Slud, E. and F. Vonta. Efficient semiparametric estimators via modified profile likelihood, *Jour. Statist. Planning and Inference*, Vol. 129, 2005, pp. 339–367.
- 12. Turnbull, B. W. The empirical distribution function with arbitrary grouped, censored and truncated data, *Journal of the Royal Statistical Society*, Vol. 38, 1976, pp. 290–295.
- 13. Van der Vaart, A. W. and J. A. Wellner. Weak Convergence and Empirical Processes (with applications to statistics). Springer Series in Statistics. Springer, 2000.
- 14. Van der Vaart, A. W. Asymptotic Statistics. Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 1998.
- 15. Vonta F. Efficient estimation in a non-proportional hazards model in survival analysis, *Scand. Journal of Statist.*, Vol. 23, 1996, pp. 49–61.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 50–56 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# UPON CONTROLLING SEVERAL BUILDING PROJECTS IN A TWO-LEVEL CONSTRUCTION SYSTEM

**D.** Golenko-Ginzburg<sup>1,2</sup>, Z. Laslo<sup>3</sup>

<sup>1</sup>Department of Industrial Engineering and Management, Ariel University Center (AUC) of Samaria Ariel, 44837, Israel <sup>2</sup>Department of Industrial Engineering and Management, Ben-Gurion University of the Negev Beer-Sheva, 84105, Israel E-Mail: dimitri@bgu.ac.il

<sup>3</sup>The Department of Industrial Engineering and Management, SCE- Shamoon College of Engineering Beer-Sheva, 84100, Israel E-mail: zohar@sce.ac.il

A two-level construction system is considered to be composed of several different building projects  $U_i$ ,  $1 \le i \le n$ , at the lower level and a control device at the upper one. The upper system's level is required to produce a given target amount V by a given due date D subject to a chance constraint, i.e. the least permissible probability p of meeting the target on time is pregiven. Each building

project  $U_i$  has several possible speeds  $v_{i1}$ ,  $v_{i2}$ , ...,  $v_{im}$ , which are subject to random disturbances. The project's output can be measured only at preset inspection (control) points. The target amount is gauged by a single measure, e.g. in square meters, and may be rescheduled among the projects. For each unit, the average costs per time unit for each project and the average cost of performing a single inspection at a control point to observe the actual output at that point are given.

We present a two-level on-line control model under random disturbances, which centres on minimizing the system's expenses subject to the chance constraint. The suggested two-level heuristic algorithm is based on rescheduling the overall target among the projects both at t = 0, when the system starts functioning, and at each emergency point, when it is anticipated that a certain project is unable to meet its local target on time subject to a chance constraint. At any emergency point t the remaining system's target  $V_t$  is rescheduled among the projects; thus, new local targets  $V_{it}$ ,  $1 \le i \le n$ ,  $\sum_i V_{it} = V_t$ , are determined. New local chance constraint values  $p_{it}$  are determined too. Those values enable the system to meet its overall target at the due date subject to the pregiven chance constraint p.

Keywords: production speed, cost-optimisation, target amount reassignment, chance constraint, inspection point

# 1. Introduction

In recent years the problem associated with developing multilevel on-line production control models under random disturbances for flexible manufacturing systems has been discussed in the literature [1–13, 15, 16]. Most of those investigations deal with not fully automated plants of 'man-machine' type where the output cannot be measured continuously on-line, but only at preset control points. The main idea of the interaction problems between different levels in hierarchical control systems is based on the conception of emergency. By using the idea that hierarchical levels can interact only in special situations, the so-called emergency points, one can decompose a general and complex multi-level problem of optimal production control into a sequence of one-level problems.

Two different optimisation cases are usually considered:

- 1. Case with a conflicting two-criteria objective, namely, to maximize the probability of completing the production on the due date, and to minimize the number of control points; but the first criterion is dominant.
- 2. The objective is to maximize the expected net profit.

Note that most of the papers outlined above do not implement a chance constraint in the on-line production control model. In our opinion, minimizing the system's expenses to meet the target on time, i.e. at a given due date, is not to be the only goal in the course of the long-term cooperation with various customers. To honour the company's good name, an additional requirement has to be inserted in the model: the production system has to meet its due date on time with a pregiven confidence probability. Thus, a chance constraint has to be implemented in the control model.

A two-level construction system is considered to be composed of several different building projects  $U_i$ ,  $1 \le i \le n$ , at the lower level and a control device at the upper one. The upper system's level is required to build a given target amount V by a given due date D subject to a chance constraint, i.e. the least permissible probability p of meeting the target on time is pregiven. Each building project  $U_i$ 

has several possible speeds  $v_{i1}$ ,  $v_{i2}$ , ...,  $v_{im}$ , which are subject to random disturbances. The project's output can be measured only at preset inspection (control) points. The target amount is gauged by a single measure, e.g. in square meters, and may be rescheduled among the projects. For each project, the average building costs per time unit for each speed and the average cost of performing a single inspection at a control point to observe the actual output at that point are given.

Golenko-Ginzburg et al. [7] have developed a cost-optimisation on-line control model which for a single project determines both control points and speeds to be introduced at those points, in order to minimize the project's expenses within the planning horizon, subject to the chance constraint. We present a two-level on-line control model under random disturbances, which centres on minimizing the system's expenses subject to the chance constraint. The suggested two-level heuristic algorithm is based on rescheduling the system's target among the production units both at t = 0, when the system starts functioning, and at each emergency point, when it is anticipated that a certain project is unable to meet its local target on time subject to a chance constraint. At any emergency point t the remaining system's target  $V_t$  is rescheduled among the projects; thus, new local targets  $V_{it}$ ,  $1 \le i \le n$ ,  $\sum_i V_{it} = V_t$ , are determined. New local chance constraint values  $p_{it}$  are determined too. Those values enable the system to meet its overall target at the due date subject to the pregiven chance constraint p.

After reassigning to each project  $U_i$  its new target  $V_{it}$  and the chance constraint value  $p_{it}$ , the projects first work independently and are controlled separately. At each k-th control point  $t_{ik}$  of project  $U_i$ , given the actual amount already produced, decision-making centres on determining both the next control point  $t_{i,k+1}$  and the index j of the new speed  $v_{ij}$  to proceed with up to that point,  $1 \le j \le m$ . The on-line control for each project proceeds either until the next emergency point, or until the due date D.

Rescheduling the remaining system's target amount  $V_t$  among the projects is carried out by using heuristic procedures. Determining chance constraint values  $p_{it}$  is carried out by using a cyclic coordinate descent method in combination with a two-level simulation model.

#### 2. Notation

Let us introduce the following terms:

S	<sup>-</sup> the two-level construction system composed of $n$ building projects $U_i$ ,
D	$1 \le i \le n$ ;
D	- the due date (pregiven),
$D_t$	<sup>-</sup> the length of the remaining planning horizon at moment t, $D_t = D - t$ ;
F	- the actual moment the target amount is completed (a random value);
р	- the chance constraint, i.e. the minimal permissible confidence probability of accomplishing the system's plan on time (pregiven);
$p_{it}$	$$ the chance constraint value for each project $U_i$ determined at the emergency
	moment $t \ge 0$ , $1 \le i \le n$ (to be determined as an optimized variable);
S <sub>ik</sub>	<sup>-</sup> the index of the speed chosen by the decision-maker at the control point $t_{ik}$ ;
$t_{ik}$	<sup>-</sup> the k -th inspection moment (control point) of project $U_i$ , $k = 0, 1,, N_i$ ;
$t_q^{em}$	<sup>-</sup> the q-th emergency moment at the system level, $l \le q \le N_{em}$ (a random value);
$N_i$	<sup>-</sup> the number of inspection moments for each project $U_i$ ;
N <sub>em</sub>	- the number of emergency moments (a random value);

V <sub>ij</sub>	-	the <i>j</i> -th speed of project $U_i$ to construct its target, $l \le j \le m$ (a random
		value with pregiven density function $f_{ij}(v)$ ;
$\overline{v}_{ij}$	-	the average of speed $v_{ij}$ . It is assumed that for each project $U_i$ speeds $v_{i1}$ ,
		$v_{i2}$ ,, $v_{im}$ are sorted in ascending order of their average values and are
		independent of t. Thus, value $\overline{v}_{im}$ is the maximal average speed for project $U_i$ ;
V	-	the pregiven system target (planned program) gauged by a single measure (target amount);
$V^{f}\left(t\right) = \sum_{i=1}^{n} V_{i}^{f}\left(t\right)$	-	the actual system's output observed at moment $t$ (a random value);
$V_{it}$	-	the target amount assigned to project $U_i$ at the emergency point $t$ (to be
		determined); note that $\sum_{i} V_{it} = V_t$ ;
$V_i^f(t)$	-	the actual output of project $U_i$ observed at moment $t$ , $0 \le t \le D$ ;
		$V_i^f(0) = 0$ (a random value);
$V_t$	-	the system's remaining target amount at moment t, $V_0 = V$ ;
$W_{p}\left[V_{i}^{f}(t), V_{it}, j\right]$	-	the $p$ -quantile of the moment target amount $V_{it}$ will be completed on conditions
-		that: (a) speed $v_{ij}$ is introduced for project $U_i$ at moment $t$ and will be used
		throughout, and (b) the actual observed output of project $U_i$ at moment $t$ is
		$V_i^f(t);$
m	-	the number of possible speeds (common to all projects);
d	-	the minimal time span between two consecutive control points $t_{ik}$ and $t_{i,k+1}$
h	_	(pregiven); equal for all projects;
$n_i$	_	the search step for determining optimal values $p_{it}$ ;
	-	the minimal value of the closeness of inspection moment $T_{ik}$ to the due date
<i>a.</i> .	-	D (pregiven and equal for all projects); lower bound of random speed $V_{in}$ (pregiven):
$a_{ij}$	-	upper bound of random speed $V_{ij}$ (pregiven);
C	_	the total operational costs, penalties and charges accumulated for the system in
C		the course of accomplishing the target amount (a random value);
$C_{em}$	-	the average cost of rescheduling the remaining target amount $V_t$ among
		projects $U_i$ by the system at a routine emergency moment $t \ge 0$ ;
$C_{ij}$	-	the average processing cost per time unit of speed $v_{ij}$ , $1 \le i \le n$ , $1 \le j \le m$
		(pregiven); note that for a fixed $i$ relation $j_1 \le j_2$ results in $C_{ij_1} < C_{ij_2}$ ;
C <sub>ins</sub>	-	the average cost of performing a single inspection of a project (pregiven, equal for all project):
$C^{f}(t)$	-	the actual accumulated processing and inspection costs calculated at moment
		t for project $U_i$ , $0 \le t \le D$ , $1 \le i \le n$ , $C_i^f(0) = 0$ ;
$C^{*}$	-	the penalty paid to the customer by the system for not accomplishing the target
C**	_	amount on time, i.e. when $F > D$ (a single payment, pregiven); the penalty cost for each time unit of dalay $F = D$ (pregiven);
C***	_	storage charges per time unit for the target amount's completion before the due
L	-	date (pregiven).

# 3. The Control Model

A two-level control model is suggested where each level faces a stochastic optimisation problem.

#### 3.1. The Problem at the System Level (Problem A)

At each emergency point  $t = t_q^{em}$ ,  $1 \le q \le N_{em}$ ,  $t_1^{em} = 0$ , determine local production plans  $V_{it}$ ,  $1 \le i \le n$ , together with local chance constraints  $p_{it}$ , in order to minimize the expected total expenses

$$\min_{\{V_{ix}, p_{ix}\}} \overline{C}$$
(1)

subject to the chance constraint

$$Pr\left\{V^{f}(D) \ge V\right\} \ge p.$$
<sup>(2)</sup>

Note that random value C satisfies

$$C = \sum_{i=1}^{n} \sum_{k=0}^{N_{i}-1} \left[ C_{is_{ik}} \left( t_{i,k+1} - t_{ik} \right) \right] + \sum_{i=1}^{n} \left( N_{i} - 1 \right) C_{ins} + N_{em} C_{em} + \left[ C^{*} + C^{**} (F - D) \right] \delta + C^{***} (D - F) (1 - \delta),$$
(3)

where

$$\delta = \begin{cases} 1 & \text{if } F > D \\ 0 & \text{otherwise,} \end{cases}$$
(4)

and values  $\{s_{ik}\}$  and  $\{t_{ik}\}$  are obtained by solving Problem B at the project level.

Values  $\{V_{it}\}\$  at each emergency point t, including t = 0, are determined according to a widely used heuristic procedure [4, 12], namely

$$V_{it} = V_t \frac{\overline{V}_{im}}{\sum_{i=1}^{n} \overline{V}_{im}} ,$$
(5)

where  $\overline{v}_{im}$  is the maximal construction speed which can be introduced for the building project  $U_i$ .

As to values  $\{p_{it}\}$ , they are determined by using one of the classical search procedures for optimising a multi-dimensional non-linear function, e.g. a cyclic coordinate descent algorithm [3–5, 14]. The search procedure is carried out via simulation, by undertaking numerous realizations of a simulation model at the lower level in order to obtain representative statistics. The simulation model represents the process of manufacturing for several building projects  $U_i$  with input values  $\{V_{it}\}$  and  $\{p_{it}\}$ , between two adjacent emergency points  $t_q^{em}$  and  $t_{q+1}^{em}$ . In the case of a routine emergency call the problem at the section level is resolved, new values  $\{V_{it}\}$  and  $\{p_{it}\}$  are determined, and the manufacturing process proceeds at the lower level, for each project  $U_i$  independently.

# 3.2. The Problem at the Project Level (Problem B)

The cost-optimization control model for a single building project has been formulated in [2, 7]. We have modified that problem for the case of several projects with additional cost parameters  $C_{em}$ ,  $C^*$ ,  $C^{**}$  and  $C^{***}$ .

For the case of an independent project  $U_i$ , given the input values  $V_{it}$ ,  $p_{it}$ , d,  $\Delta$  and  $\overline{v}_{ij}$ ,  $l \leq j \leq m$ , the problem is to determine both control points  $\{t_{ik}\}$  and building speeds  $\{v_{is_{ik}}\}$  to minimize the construction expenses

$$J = \underset{\{t_{ik}, v_{is_{ik}}\}}{\min} \left\{ \sum_{k=0}^{N_i - 1} \left[ C_{is_{ik}} \left( t_{i,k+1} - t_{ik} \right) \right] + N_i C_{ins} \right\}$$
(6)

subject to

$$Pr\left\{V_{i}^{f}\left(D\right)\geq V_{ii}\right\}\geq p_{ii},$$
(7)

 $t_{i0} = t , (8)$ 

$$t_{iN_i} = \min_{T_i} \left[ T_i : Pr\{V_i^f(T_i) \ge V_{it}\} \right], \tag{9}$$

$$t_{i,k+l} - t_{ik} \ge d , \qquad (10)$$

$$D - t_{ik} \ge \Delta, \ 0 \le k \le N_i - l , \tag{11}$$

$$s_{ik} = j = \min_{1 \le q \le m} q \quad \forall q : W_p \Big[ V_i^f(t), V_{it}, q \Big] \le D \quad .$$

$$\tag{12}$$

Restriction (8) means that after reallocating target amounts at the routine emergency point t, the starting moment to proceed constructing, i.e. the first control point to undertake decision-making and to determine  $s_{i0}$  and  $t_{i1}$ , is t. Note that at all emergency points the remaining target amount, as well as the due date, are updated, i.e. the ordinate t = 0 is shifted to the right. Restriction (9) means that the last inspection point is the moment target amount  $V_{it}$  is reached. Restrictions (10) and (11) ensure the closeness between two consecutive control points, as well as the closeness of the routine inspection point to the due date. Restriction (12) means that the speed to be chosen at any routine control point  $t_{ik}$  must not exceed the minimal speed which guarantees meeting the deadline D on time, subject to the chance constraint (7).

The general idea of solving problems (6)–(12), which is a very complicated stochastic optimization problem, is as follows. At each control point  $t_{ik}$  decision-making centers on the assumption [7] that there is not more than one additional control point before the due date. Two speeds have to be chosen at point  $t_{ik}$ :

- 1. Speed  $v_{ij_l}$ ,  $j_l = s_{ik}$ , which has to be actually introduced at point  $t_{ik}$  up to the next control point  $t_{i,k+l}$ .
- 2. Speed  $v_{ij_2}$ ,  $j_2 = s_{i,k+1}$ , which is forecast to be introduced at control point  $t_{i,k+1}$  within the period  $\begin{bmatrix} t_{i,k+1}, D \end{bmatrix}$ .

Thus,  $j_1$  is determined in accordance to Eq. (12) and  $j_2$  is determined by honoring chance constraint (7). In [7] at each routine control point  $t_{ik}$ , all possible couples are singled out. The couple, which delivers the minimum of forecasted manufacturing and control expenses, has to be chosen. Since couple  $(j_1, j_2)$ , together with the inspected value  $V_i^f(t_{ik})$  and values D and  $V_{it}$ , fully determines the next control point  $t_{i,k+1}$ , speed  $v_{ij_1}$  is introduced within the period  $[t_{ik}, t_{i,k+1}]$ . At moment  $t_{i,k+1}$  decision-making has to be carried out anew.

# 4. The General Idea of the Two-Level Heuristic Algorithm

The general idea of the suggested heuristic algorithm is as follows: at each routine emergency point  $t_q^{em}$ ,  $q = 0, 1, ..., N_{em}$ , decision-making centres on minimizing the *future costs* from point  $t_q^{em}$  until F, including the penalty and the storage costs. The costs representing the past (interval  $[0, t_q^{em}]$ ) are not relevant for this on-line control problem, and there is no need to remember the past decision [9]. The only relevant information to be stored is  $t_q^{em}$  and  $V_i^f(t_q^{em})$ . Thus, decision-making at the system level is carried out only at emergency points  $t_q^{em}$  including the moment t = 0 the system starts constructing.

Decision-making at the system level at each routine emergency moment  $t = t_q^{em}$  centers on determining both new chance constraint values  $\{p_{it}\}$  and new target amounts  $V_{it}$  for the remaining planning horizon [t, D]. Values  $\{p_{it}\}$  are obtained via simulation, by a combination of a search algorithm and an on-line one-level control algorithm for several projects. The latter work independently and are controlled separately at inspection points. It is generally assumed that at the beginning of the work all the available resources, i.e., the building teams, are previously allocated among the projects. Those resources remain unchanged within the planning horizon, i.e. no resource reallocation is performed. Thus, the corresponding construction speeds  $v_{ii}$  for each project  $U_i$  remain unchanged too.

If for a certain project  $U_i$  at a routine inspection point  $t_{ik}$  it is anticipated that the project cannot meet its target  $V_{it}$  on time subject to the previously determined chance constraint  $p_{it}$ , an emergency is then called, and decision-making is affected at the system level. The remaining target  $V_t$  at  $t = t_{ik}$ , together with the remaining time  $D_{t_{ik}} = D - t_{ik}$ , is then updated. New quasi-optimal values  $\{p_{it}\}$ ,  $t = t_{ik}$ , together with new target amounts  $\{V_{it}\}$ , are then determined. The newly corrected plan is assigned to all building projects, and the construction process proceeds further, until either the new emergency point or until the moment the target amount is completed. Thus, decision-making at the system level centers on numerous recalculations of the system's plan subject to the chance constraint. This is carried out by using a forecasting simulation model with input values  $\{V_{it}, p_{it}\}$ ,  $t = t_{ik}$ . The matrix  $Z = \{V_{it}, p_{it}\}$  which delivers the minimum of total accumulated costs subject to the chance constraint p, is taken as the *optimal corrected plan*. Afterwards, that corrected plan is passed to the projects, and on-line decision-making is carried out at the project level.

#### **Conclusions and Future Research**

The following conclusions can be drawn from the study:

- 1. The two-level control model under consideration is a further development in the area of production control. The suggested control algorithm enables meeting the target amount on time subject to a chance constraint for two-level man-machine production systems under random disturbances.
- 2. The developed control model can be used for various semi-automated production systems [1, 9] under random disturbances where the outputs (target amounts) can be measured only at pregiven control points and are gauged by a single measure (money terms, square or cubic meters, percentage completion of the project, etc.). When a building construction company performs a building project, units are teams, which can work with several alternative speeds, while the output of a team is usually measured in percentage completion of the project.
- 3. The system's target amounts are transferable and may be rescheduled among the production units.
- 4. For all previously developed multilevel semi-automated production systems under random disturbances the risk average principle has been implemented in the control model at the lower level [1, 3–6, 8–13]. Those models do not deal with chance constraints. The model under consideration is based on another principle, namely, the chance constraint principle [2, 7] which is very effective for cost objectives.
- 5. Future research may be undertaken to develop three-level control models, e.g. factories, etc., with chance constraints at each hierarchical level.

# References

- 1. Golenko-Ginzburg, D. A two-level production control model with target amount rescheduling, *J. Oper. Res. Soc.*, Vol. 41, No 11, 1990, pp. 1021–1028.
- Golenko-Ginzburg, D. and A. Danoch. Cost-optimization on-line production control model. In: *Proceedings* of the 8<sup>th</sup> International MIRCE Symposium on System Operational Effectiveness, Exeter, UK, 7–9 December, 1998, pp. 315–323.
- 3. Golenko-Ginzburg, D. and V. Kats. A generalized control model for man-machine production systems with disturbances, *Comput. Ind. Eng.*, Vol. 32, No 2, 1997, pp. 399–417.
- 4. Golenko-Ginzburg, D. and V. Kats. A generalized production control reallocation model. In: *Intelligent Scheduling of Robots and Flexible Manufacturing Systems / E. Levner (Ed.).* Holon, Israel, 1997, pp. 139–156.
- Golenko-Ginzburg, D. and V. Kats. A three-level factory control model. In: Proceedings of the Israel-Germany Bi-National Conference on Computer Integrated Extended Manufacturing Enterprise, 28–29 February, 1996.
- Golenko-Ginzburg, D., Kats, V. A generalized control model for man-machine production systems with disturbances. In: *Proceedings of the 13<sup>th</sup> International Conference on Production Research*, Jerusalem, 6–10 August, 1995, pp. 590–592.
- 7. Golenko-Ginzburg, D., Gonik, A., Papic, L. Developing cost-optimization production control model via simulation, *Mathematics and Computers in Simulation*, Vol. 49, 1999, pp. 335–351.
- 8. Laslo, Z. Activity time-cost tradeoffs under time and cost chance constraints, *Computers & Industrial Engineering*, Vol. 44, No 3, 2003, pp. 365–384.
- Laslo, Z. and D. Golenko-Ginzburg. Resource models for several flexible transportation systems under random disturbances with central warehouses, *Computers Modelling and New Technologies*, Vol. 7, No 2, 2003, pp. 48–50.
- 10. Laslo, Z., Golenko-Ginzburg, D. and A. Gonik. Alternative stochastic network projects with renewable resources, *Computers Modelling and New Technologies*, Vol. 9, No 1, 2005, pp. 40–46.
- 11. Laslo, Z., Gurevich, G. Minimal budget for activities chain with chance constrained lead-time, *International Journal of Production Economics*, Vol. 171, No 1, 2007, pp. 164–172.
- 12. Laslo, Z., Keren, B., Ilani, H. Minimizing task completion time with the execution set method, *European Journal of Operational Research*, Vol. 187, pp. 1513–1519.
- 13. Keren, B. A dynamic algorithm which integrates a static solution with the 'EVPI' principle to determine inspection points and production speeds: Master Thesis. Department of Industrial Engineering and Management, Ben-Gurion University of the Negev. Beer-Sheva, Israel, 1989.
- 14. Luenberger, D. G. Linear and Nonlinear Programming, 2<sup>nd</sup> edition. MA, USA: Addison-Wesley, 1989.
- 15. Sethi, S. P., Zhang, Q. *Hierarchical Decision-making in Stochastic Manufacturing Systems*. Boston, Cambridge, MA: Birk-hauser, 1994.
- Sinuany-Stern, Z., Golenko-Ginzburg, D., Keren, B. A static adaptive approach for solving a dynamic problem of production control. In: *Proceedings of the 3<sup>rd</sup> European Simulation Congress / D. Murray-Smith, J. Stephenson, R. N. Zobel (Eds.).* Edinburgh, Scotland, 1989, pp. 509–511.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 57–64 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# PERFORMANCE EVALUATION OF TEAMWORK REFLEXIVE ANALYSIS

G. I. Petkov<sup>1</sup>, V. A. Iliev<sup>2</sup>

<sup>1</sup>Technical University of Sofia Kliment Ohridsky Boulevard 8, 1797 Sofia, Bulgaria E-mail: gip@tu-sofia.bg

<sup>2</sup>International Institute for Social Communication and Behaviour 5800 Pleven, Bulgaria E-mail: aasenova@abv.bg

The paper presents the development and supplementation of an advanced Human Reliability Analysis technique – Performance Evaluation of Teamwork method with approaches and mechanisms of reflexive analysis for prevention of erroneous *communication and decision-making*.

Keywords: reflection, reflexive analysis, typing, scenario analysis, styling, communicative barriers, Human Reliability Assessment (HRA), risk, human erroneous actions, Performance Evaluation of Teamwork (PET) method

# **1. Introduction**

The efficient operators' teamwork, as group interaction, emanates some fundamental characteristics – integration, coordination, synchronization, communication and leadership. The crew of operators needs successful integration of individual thoughts and actions. Every operator has a specific and unique role and his correct performance contributes to the crew's success. It means that the causes for crew failure could be due not only to inability and unreliability of each individual crew member but also to the crew error in coordination and synchronization of the operator's individual contribution to the common work.

The group mental process becomes especially critical for the whole crew performance in risky, complex and dynamic situations as accidents and incidents when the crew is under stress and the urgency of tasks is frequently prerequisite for inclusion of other variables and initiators for erroneous actions. In such work conditions the situation is very dynamic and the crew is threatened by unforeseen circumstances and circumventions (e.g., equipment failures, operator's violations and errors). The operators are exposed to an increasing information stream and the workload also increases. The way out is good communication that allows for correct tasks distribution, coordination and manipulations synchronization. It also makes possible the addition of cognition individual processes and enables the group leader to make objectively necessary decisions in the given situation.

From psychological point of view communication is a complicated and multidimensional space of a new "psychic reality" of contact. Existence of alternatives in the development of that reality is a premise for risk occurrence. Risk is related to the situation, individual, teamwork and social contexts development that affects considerably the relation objective-result.

Considerable part of all operators' erroneous actions is based on their communication. One of the most effective tools for social cognitive study of this communicative basis is the method of reflexive analysis [1].

The paper presents the issues of dealing with cognition and communication aspects of operators' performance. They are evaluated on the base of macroscopic "second-by-second" context model of cognition and communication processes. The human action context is represented as a field of interaction between humans, machine, technology, organization and society. All together they form the fieldwork and this context influences dynamically their configuration. Some useful achievements of mathematical psychology are applied to context quantification and dynamic reconfiguration of cognition and communication processes modelling. A simplified probabilistic approach to dynamic quantification of operators' performance and communication errors identification is demonstrated by the Performance Evaluation of Teamwork (PET) method. This approach is illustrated by retrospective human reliability analysis (HRA) of the accidents based on scenario chronological archives of the FSS-1000 (full-scope simulator of WWER-1000) of Kozloduy nuclear power plant (NPP) obtained during the regular practical training of crews.

# 2. Operators' Erroneous Action Classification and Reflexive Analysis

The operator's work contains actions based on different hybrid interactions in local "man-machine", "man-technology", "man-organization", "man-society" and "man-man" systems. Communication between people for the global hybrid system is especially very important feature because it constitutes, organizes and configures all global system interactions, determines different forms of functioning and guarantees functional redundancy and diversity.

Operator's erroneous actions can be relatively insignificant, with local influence on operator's work but they also can be crucial in case of serious accidents and emergency situations. It is a matter of concern of the safety analyses of organization, technology and performance, whether the crew communication causes errors or makes the operator's work secure and reliable.

The operator's errors strongly depend on operator's knowledge and performance manner. The most of current HRA methods distinguish between two basic operator's error modes – errors of cognition and errors of execution. Such concepts of Human Erroneous Action (HEA) were used by PET method to propound the HEA classification that is shown on Figure 1. The communicative constituents of operator's erroneous performance could be insufficient or imperfect knowledge of working matter and communication manner, incorrect communicative actions on the executive level during the teamwork performance. The classification on Figure 1 also presents another popular version of HEA distinction – between errors of omission (EOO) and errors of commission (EOC).

The operator's errors study must cover social instruments of cognition that give an advantage to the analysis of social interaction between operators. Such instruments are the methods of reflexive cognition and reflexive analysis. The reflection is especially useful in the investigation of interaction. It is a form of "feedback" and an important mechanism for getting insight into the way human psyche functioning projects the operator's active attitude to reality.



Figure 1. Human erroneous action classification

As seen on Figure 1, the subject of reflexive analysis could be all feedbacks in realization of human actions – latent errors, violations and recoveries. They appear in an iterative individual cognition form and crew communication form. These forms could be identified and distinguished experimentally by running different accident and normal scenario on the full-scope, multifunctional or computer-based simulators of NPP. This identification can be considered as the initial point of reflexive analysis application.

# 3. Reflexive Analysis

#### 3.1. Basic Features of Reflexive Cognition

A reflexive cognitive process aims to understanding the rules that individuals employ to reveal meanings externalised in their actions. The social world is a world based on meanings and images. The reflection is a way of social cognition that assigns motives, social meanings and values to people

in order to reveal motives and values that give rise to certain actions. The social cognitive process should be investigated concurrently and sequentially on the base of following mechanisms: typing, scenario analysis, styling and individual history. Some communicative barriers between operators that also cause erroneous actions could be identify and overcome by reflexive analysis.

As mentioned above, the reflection is a form of "feedback". Its cognitive meaning includes cognitive knowledge about a given object obtained by the change of positions of examination. This change of two positions, or moving to third "independent" position, allows for getting profound and comprehensive knowledge about the object. It is very important for investigation and prevention of operator's erroneous actions, especially where performance shaping factors (PSF) have complex configuration and hierarchy of casual-consecutive connections.

The most traditional illustration of reflexive knowledge is a type "I know that you know that I know". The social world is organized on the base of objective meanings. So our thoughts and conclusions about an operator's error as a kind social event depend on our subjective interpretation of the actions of other people. The reflection allows peculiar "distinguishing" of the object and provides a "look from aside". Such an individual or collective alienation from the situation creates a virtual form of reality. If we can look at the situation from aside, so we could be transformed into objective judges of what is happening.

Reflection allows to the subject of social perception to reconstruct elements of the inner world of other people. Through the creation of the image of the other's inner world operator can lead imaginary communication with another operator as a preliminary or sequent step to the real communication. By such mental modelling and transformation an operator acquires information that goes beyond the simple reflection of a separate communicative act [2].

It is necessary to choose a new point of view in order to make a reflexive look from aside to a specific situation. To choose a specific point of view it is necessary to estimate its usefulness for getting an insight into the situation, for making it clear and well understood. The more various reflexive positions the analyst (the investigator of operator's errors) possesses, the more successfully he will analyse the specific situations of erroneous actions.

In a communication situation, where operator's errors are available, there is a bilateral process of reflection between the operator, that has acted erroneously and the analyst of these erroneous actions. These new reflexive positions include: the operator as he is in reality; the operator as he imagined himself; the operator as he is imagined by the other operator; the operator as he is imagined by the investigator.

# 3.2. Reflexive System Approaches

In the particular case of reflection of the communication between two operators, the thoughts of both operators are put to the reflection in order to understand reasons for their decisions for actions. As a result, a common reflexive system is organized that generalizes and analyses the information obtained by social cognition.

Any activity can be investigated on the base of its communicative context that is transformed in the object of reflection. The communicative reflection is considered as a situational reflection of communication.

The reflexive system accumulates information of cognitive and practical nature. This information is obtained, organized and used on the base of teleological, morphological and genealogical approaches. These are approaches for primary analysis of the investigated object to reveal its basic features.

*Teleological Approach.* This approach consists of: setting one's own goal and the object's goal, investigating the goal-situation vector, evaluation of the observer's impact on the reality, evaluation of synchronous relations, analysis of the entire mutual exchange "inside" and "outside" the object.

Teleological (purpose-oriented) approach begins with the question: "Why?" It is oriented to systems that have abilities to identify their own goals (including the means to achieve them), opportunities and limits. The operators' work is a system of this kind in spite of all the diversity of behaviours.

*Morphological Approach.* It is discovering the most important quality of the object structure from the point of view of the specific analysis and object development and introducing the systematic principle for investigation of different levels of the object's existence.

Morphological approach to reflection studies is interested in the "full field of knowledge" of a given object. The methods of negation and construction are especially effective. In the end, more thorough and comprehensive knowledge about object could be achieved that is not reduced only to extreme standpoints and covers all shades and "soft" manifestations.

*Genealogical Approach.* It is an analysis of the tendencies in the family biography by using the principle of genesis of single biographical facts.

When genealogical approach is applied, the reflexive cognition is interested in the origin of social object of study. The origin is the historical biography and "genealogical" relations to the phenomena in social periphery of cognition. The problems into consideration are: "Did similar operator's errors happen in the past?", "What do they look like?", "Which phenomena are they related to and different from?", "What does usually happen during this operator's error?" etc.

The reflection cognition determines the structure of reflexive analysis [2].

#### 3.3. Reflexive Mechanisms

In reflexive analysis, the social cognitive process is developed on the base of four mechanisms: typing, scenario analysis, styling and individual history. The object is examined through the prism of these four mechanisms from the new point of view [2].

## 3.3.1 Typing

Any world interpretation is connected with our preliminary "familiarization" with it. In the subjective field of meanings the typing is connected with repetition; similarity; locality in linear causality; associative analogy of likeness. In comparison to the typical symptoms related to the object, where both sides of interpresonal communication are not coherent, we could use two modes of reflection. The first one consists in change of the point of view of one operator to the point of view of other. The second mode uses coincidence of the relevant point related to the object.

The typing is based on stereotypes and human individual experience. A rich social experience can help the analyst of operator's erroneous actions to localize more exactly the point of view to the typical symptoms of erroneous actions and cognition and execution connected to it. Additionally, it should be taken into account that such important symbolic systems as myth, tradition, authority, ritual and prejudice are powerful means for understanding of typing.

The typology of the object is based on the principle of analogies: outer, behaviour, communicative, situation and role analogy.

#### 3.3.2 Scenario Analysis

Scenarios are actions in sequence, similar in different situations. There is a variety in behaviour scenarios but in most of cases they differ only in details and some negligible symptoms. The scenarios depend both on outer conditions and the psychic qualities of a person. There are two scenario modes: pragmatic-reflexive and theoretical-reflexive.

Some examples of pragmatic-reflexive scenarios are: "there is no time for anything", "nothing comes of it again", etc. Theoretical-reflexive scenarios can be described as proactive, reactive, with open or closed end, etc.

In reflection cognition, it is appropriate to divide scenarios into outer and inner ones. The outer scenarios are imposed on a subject by outer sources, such as norms, professional rules, other subjects, etc. The inner scenarios are mental models of the sequence of actions that leads to a given result. Both types of scenarios of operator's errors must be investigated. The reflexive analysis reveals the way of refraction of objective to subjective in the operator's consciousness and to what extent he has controlled his behaviour in the specific working situation.

Scenarios related to the object are scenarios of development and behaviour.

# 3.3.3 Styling

The styling is biographically determined behaviour that reflects the individual predisposition to specific actions in specific situations. The individual style is the usual way a person solves of typical problem.

In conformity with one's individual style the operator uses the maximum compensatory mechanisms in order to apply all necessary qualities to achieve a given purpose. The individual style of operator is a very important feature for the analysis of causes for operator's error. The modification of individual style is not impossible but it is a very difficult task. That is why the management of an operators' crew should preliminary take measures for synchronization of individual styles of the crew members. The individual style of communication can change for different people in different features. For example, some operators have a formal style of communication and others show an intimate, friendly style. It should be noted that the style of communication can reflect directly the causes of erroneous actions.

# 3.3.4 Individual History

The individual history includes specific events in a person's life and his own attitude to them. Every human action is biographically determined. The individual history is often perceived as a "personal lifetime". Its reflexive analysis allows finding a person's driving forces, worldly perspective, and the often used ways of behaviour. As a result, a frequency register of the most important behaviors and the most typical responses of the operator are obtained. The operators that make errors are different but everybody is relatively identical to one's individual history. Because of that we can prognosticate an operator's resilience against erroneous actions in hybrid system communication based on a preliminary study of his individual history.

# 4. Evaluation of Crew Communication in Main Control Room by the PET Method

The Performance Evaluation of Teamwork method gives opportunity to take into account dynamical differences in operators' achieved knowledge during their cognitive process and the arising need of natural communication between crew members [3]. It consists of two reliability models of processes of cognition and communication, represented as directed graphs. They are based on quantification of individual context of each operator in situation by successively applying of methods "Violation of Objective Krebs" and "Combinatorial Context Model". A quantified individual context probability is used to obtain a cognitive error probability (CEP) of the success of individual cognition for each operator and the success of decision-making for operator's crew.

A probability of natural communication (PNC) arising between two operators of the crew is evaluated as difference of individual CEP of these operators.

Because of the complexity, dynamics and risk of processes on NPP it is obligatory to apply an initiated communication not only to the crew working in the main control room (MCR) but also to the conversations between the related operators outside the MCR, such as supervisors, dispatchers, etc. The workplace instructions, technological regulations, working and emergency procedures describe not only the licenses, duties and subordination of each operator but also specific orders and reports (communication) that operator must perform in a given situation. This initiated communication could be also modelled and its impact on the crew error probability can be evaluated by the PET method.

A probability of initiated communication (PIC) between two operators of the crew is accepted as a complement probability to the PNC (PNC + PIC = 1).

The study issues concerning natural and initiated communication impact on three MCR crews reliability are shown briefly below. Each of the investigated operators' crews ran twice the accident scenario 'Steam Generator Tube Rupture'' (SGTR) during their regular training on FSS-1000 of Kozloduy NPP with WWER-1000 (Russian pressurized water reactor).

## 4.1. Special Features of Scenario and Its PET Model

The study data of the PET model are obtained from the FSS-1000 training program, operation and accident procedures of the Kozloduy NPP, unit 6. The available computer means of FSS-1000 made possible the on-line data mining for registered events, controlled technological parameters and monitored operator's and instructor's actions. The possibilities of digital audio-visual system (AVS) were used for: operator's cognitive processes modelling, prognosis of probable performance shaping factors (PSF) influence to the crew operator's actions, communication and decision-making and the evaluation of their effect. Because of the lack of individual microphones to record the conversations between operators during scenario running the timing of the unit supervisor orders and operator's report to them were filled in preliminary prepared written forms by the FSS-1000 instructor monitored training classes by the AVS.

The PET communication reliability model used for MCR operator crews of NPP with WWER-1000 is shown on Figure 2.



Figure 2. Main Control Room Crew Communication Reliability Model of NPP with WWER-1000

The training class description and accident timing were presented partially in [4] and fully in [5]. Two crews and related operators (6 investigated archives in total) have performed repeated erroneous actions (mistakes or violations) that have changed the accident course and need special attention with respect to the used equipment and procedures. The erroneous action is "Unclosed valve in accident gas freeing system" (YR) that is applied for depressurising of primary circuit. It leads to the membrane breaking of bubbler tank and increasing of pressure and humidity at the containment. The admission of two leaks from the primary circuit (from steam generator and primary valves) is an accident beyond design basis because two design accident events had coincided. The results for efficiency of communication for the crew reliability (the negative difference between crew CEP) are shown on Figure 3 and 4 below for the SGTR scenario with violation and without violation respectively.

The mental processes of cognition, communication and decision-making are iterative processes. With regard to the model simplicity and conservatism one-step iteration in both digraph reliability models is used. They include Rasmussen's Step-Ladder (SLM) model of individual cognition [3], model of the crew's mutual communication and decision-making with the assumption that the supervisor (decision-maker) is absolutely reliable if he has the entire information about the situation. This assumption should be made because a PET model for decision-making is not available yet.



The List of communications of crew "17a" and their durations in scenario SGTR are presented in Table 1 below.

Figure 3. Communication Efficiency of 17a Crew with violation





Table 1. List of communications of "17a" crew in scenario "Steam Generator Tube Rupture"

Ν	Duration, s	Content
1	132–140	Supervisor orders to SRO/RO to scram reactor because of leak ">501/h"
2	176–286	Supervisor orders to SRO to turn off the first main coolant pump (MCP1)
3	202–284	Supervisor orders to OFWP/TO to isolate first stem generator (SG1) and turn off 1-of-3 turbo-feed-water pump.
4	212–382	Supervisor orders SRO/OR to put on regulator YP04 in automatic mode and cooling with rate of 20°C/h"
5	382-807	Supervisor orders to OT to start emergency cooling with rate of 60°C/h through steam dump system to the condenser up to 52kgf/cm <sup>2</sup>
6	982-1101	Supervisor orders the TO to close first quick closing valve (TX50S06) and its bypasses (TX50S23,24) in 70–75 kgf/cm <sup>2</sup>

#### 4.2. PET Communication Study Results

The study demonstrates indisputably that the teamwork decision-making has advantages over the individual decision-making. It is especially efficient when mistake or violation occurs (see 17a vs. 18a on Figures 3 and 4). In these cases the role of communication for decision-making support enhances.

The comparison of initiated and natural communication (PIC and PNC) shown on Figure 3 demonstrates that the PIC has very positive influence on successive cognition and decision-making. However, the PET analysis makes it evident that PIC could be in time or late and the positive influence reduces for the latter case.

The efficiency of communication for the crew reliability is defined as  $\Delta CEP=CEP_{Supervisor}-CEP_{Crew}$ . The positive values of  $\Delta CEP$  mean that Supervisor completes the individual cognitive process and negative values mean that it is not completed. The PET communication study demonstrates the positive influence of natural and initiated communication, so the good practice requires active crew communication and synergic teamwork [6].

# 5. Reflexive PET Communication Analysis

This is an analysis of the new, externalised point of view obtained on the base of investigation of the main features of the object. It follows from the definitions of natural and initiated communication that PNC is determined in the first place by the rates of iterative individual cognitive processes of the crew members and PIC depends on operator's individual willing and interpersonal relations. These factors contribute mostly for the active communication between the operators in a crew.

As mentioned above, this is an initial point for reflexive communication analysis by the PET method. The PET evaluations for PNC and PIC present a rather rough guess, so they require qualitative interpretation based on features, approaches and mechanisms of the reflexive analysis that were described in the previous section.

As a results of the PET analysis based on four scenarios are were run by six crews two times varying one PSF at each scenario, different individual and crew CEP, PNC, PIC and communication efficiency were obtained. After that it is necessary to use reflexive analysis to explain qualitatively the variation in these communicative results and in the abilities of the operators in one crew or in different crews but performing the same job.

The variation in PET cognition and communication results are explained by qualitative position analysis and the following position modes:

- *'Fusion' Position.* It consists of identification and imitation.
- *'Way out' Position.* This position follows adverse, ambivalent and love-hate position or positions different in major qualities or in minor qualities.
- Overcoming the 'Astigmatic' Position. It is necessary to overcome the "astigmatic" position, i.e. the position of different roles, such as "the professional is prone to extreme stand-point".

However, the formal presentation of the reflexive position analysis implementation and the systematisation of the PET results based on scenario analysis must be preceded by collecting additional information about operators by using also other reflexive mechanisms as typing, styling and individual history [7].

# Acknowledgements

The research and presented paper are supported by the National Science Fund, Ministry of Education and Science of Bulgaria, project No. HS-I-308/2007 "Development of training program, data collection and performance evaluation system of students on a NPP computer simulator with WWER".

# References

- 1. Iliev, V. A. Risk and Communication. Sofia: Lege Artis, 2004. (In Bulgarian)
- 2. Iliev, V. A. Social Cognition Borders of Reflection. Sofia: Lege Artis, 2008. (In Bulgarian)
- 3. Petkov, G. Retrospective human reliability analysis by performance evaluation of teamwork method. In: *Proceedings of the PSAM6 Conference, San Juan, USA, 23–28 June, 2002.* San Juan, USA, 2002.
- 4. Pekov, M., Petkov, G. Evaluation of faster human-machine interface on main control room operators' crew reliability in "Steam Generator Tube Rupture" accident. In: *Proceedings of EMF'06 conference, Varna, 18–20 September, 2006.* Varna, 2006, pp. 199–212. (In Bulgarian)
- 5. Pekov, M. Quantitative assessment of human reliability by the PET method: Analysis of operator's actions during accident "Steam generator tube rupture of WWER-1000" on the base of FSS-1000 archive data: Thesis for B.Sc. degree, TU-Sofia, July 2006. Sofia, 2006. 70 p. (In Bulgarian)
- 6. Petkov, G., Iliev, V. Reflexive analysis of operator's erroneous actions. In: *Proceedings of the SMRLO'10 Conference, Beer Sheva, Israel, 8–11 February, 2010.* Beer Sheva, Israel, 2010.
- Iliev, V. and G. Petkov. Confidential communication for reliable information and knowledge exchange. In: *Proceedings of the PSAM10 Conference, Seattle, USA, 7–11 June, 2010.* Seattle, USA, 2010. (In press)

Received on the 21<sup>st</sup> of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 65–72 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

# THE HOTELLING'S METRIC AS A CLUSTER STABILITY MEASURE

Z. Volkovich, Z. Barzily, D. Toledano-Kitai, R. Avros

Software Engineering Department ORT Braude College of Engineering Karmiel, 21982, Israel E-mail: vlvolkov@braude.ac.il, zbarzily@braude.ac.il, dvora@braude.ac.il, r\_ avros@braude.ac.il

The estimation of the suggested number of clusters in dataset is an ill posed problem of essential relevance in cluster analysis. A group (cluster) is characterized by a relatively high similarity among its elements in addition to a relatively low similarity to elements of other groups. High stability in partitions, obtained from the same data source, is logically classified as a high consistency of the clustering process. Thus, the number of clusters that maximizes cluster stability can serve as an estimator for the "true" number of clusters. In the current paper we consider a probabilistic approach to this problem resting upon the Gaussian clusters model. We claim that sequences of clustered samples can be interpreted as Gaussian distributed i.i.d. samples drawn from the same source, if the number of clusters is chosen correctly. The samples closeness, within the clusters, can be measured by means of the *p*-values, calculated for the appropriate Hotelling's *T*-square statistic. Data outliers and clustering algorithm's shortcomings can yield small *p*-values. However, we presume that their empirical distribution is least concentrated at the origin, if the number of clusters is chosen correctly. Our procedure can roughly be described as a generation of such a distribution with a sequential application of a *concentration test*.

Keywords: clustering, Hotelling's criteria, cluster validation

#### **1. Introduction**

Cluster analysis is an essential tool, in machine learning, commonly employed in order to identify meaningful groups, named clusters. Items belonging to the same cluster are expected to be more similar one to another in comparison with elements belonging to different clusters. Most widespread iterative clustering algorithms, such as *k*-means, *EM* and *k*-medoids, are typically carried out in three steps. The first one is the initialisation step. It is intended to set an initial partition. In the second step the data set is partitioned to the "best" possible clusters. Here, elements are assigned to clusters so that some predefined function is optimised. The third step compares the attained partition to the previous one. If the difference is less than a predefined stopping threshold the algorithm stops, otherwise it returns to the second step. The partitioning phase is usually expressed by assigning a label to each element. This label identifies the cluster to which the element belongs. In general, the elements of the labelling set do not have specific meanings thus, they can be permuted from one instance to another.

Another clustering method-type is designed to yield the optimal ("true") number of clusters. The current paper is addressed to handle this problem, which is recognized as an "ill posed" one [20], [14]. For example, an answer here can depend on the scale in which the data is measured (see, for example [6]). Many approaches have been offered. Hitherto, none of them has been accepted as superior.

From the geometrical point of view cluster validation has been considered in [11], [19] (*C*-index), [4], [18], [21], [33], [13], [26] and [34] (*the Gap Statistic method*). Stability based methods measure the cluster variability under repeated invocations of a clustering algorithm on samples drawn from the same data source. Low variability in partitions is interpreted as high consistency in the obtained result [7]. So, in papers [24], [2], [3] the stability is evaluated by the fraction of times that a pair of elements maintain the same membership under a rerun of the algorithm. In [27] a stability function uses the Loevinger's measure of isolation. The prediction-based resampling method Clest [10] employs, de-facto, external partition correlation indexes as a stability magnitude. In papers [29] and [23] a theoretical justification, of the cluster validation problem, by means of the stability concept has been discussed.

Nonparametric estimation of the underlying density provides another methodology to the determination of the "true" number of clusters. This approach assumes that the clusters correspond to density peaks. Thus, the clustering procedure appoints each item to a "domain of attraction" of the density modes. Evidently, Wishart [37] offers to look for the modes, as the first step in order to discover a cluster structure. He suggested that clustering methods must be able to expose and "resolve distinct data modes,

independently of their shape and variance". Based on this idea Hartigan ([16] Section 11 and [17]) introduced the notion of "high density clusters". The clusters amount has been defined as the quantity of disjoint areas where densities go beyond a given threshold. So, clusters are presented by islands of "high" density in a sea of "low" density. This concept has been later employed in numerous papers (see, for example [8], [9] and [32]).

A methodology based on the goodness of fit test procedure has been considered in [28], [15]. Other approaches relying on the goodness of fit tests have been offered in [35]. Here, the etalon cluster distributions are constructed using a model intended to represent a mixing of samples within the clusters. In [1], such a model has been provided by means of an asymptotically normally distributed behaviour of the edges, connecting points from different samples, in a Minimal Spanning Tree built in a cluster. The paper [36] uses the binomial distribution to model the quantities of the *K*-Nearest Neighbours belonging to the own point's sample. In the works [35] and [1] partitions stability has been described by means of distances between clustered samples drawn from cluster cores versus the ones drawn from the whole population.

In the current paper we consider a probabilistic approach, to this problem, resting upon the Gaussian clusters model. Our concept proposes a statistical homogeneity, of these simulated samples, in the case of the "true" number of clusters. Given a clustering algorithm, we asses the "true" number of clusters via the stability of the whole dataset partition, constructed by the algorithm. Indeed, we assume that the clustering algorithm reflects the inner data structure, and the stability is characterized by samples "drawn" from clusters or samples occurrences in the clusters. We claim that sequences of clustered samples can be interpreted as normally distributed i.i.d. samples, if the number of clusters is chosen correctly. The samples closeness within the clusters is measured by means of the values calculated for the appropriate Hotelling's *T*-square statistic. Data outliers and clustering algorithm's shortcomings can cause small *p*-values. However, we presume that their empirical distribution is least concentrated at the origin, if the number of clusters is chosen correctly.

#### 2. The Cluster Model

The main clustered object is a data space X considered as a finite subset of the Euclidean space  $\Re^d$ .  $C_k$  denotes the set  $\{1,...,k\}$  and  $\psi_k$  is the set of all possible permutation of the set  $C_k$ . A partition (clustering)  $\Pi_k = \{\pi_1,...,\pi_k\}$  of X is a set of k non-empty clusters such that:

$$X = \bigcup_{i=1}^{k} \pi_{i}, \text{ and } \pi_{i} \cap \pi_{j} = \emptyset, i \neq j.$$

For a given partition  $\Pi_k$ , the label function  $\alpha(\Pi_k)$ :  $X \to C_k$  assigning points to the clusters is defined as  $\alpha(\Pi_k)(x) = i$ , when  $x \in \pi_i$ , i = 1, ..., k, so  $\pi_i = \{x \in X \mid \alpha(\Pi_k)(x) = i\}$ .

Each partition of the set X provides a mixed decomposition of the underlying distribution  $P_X$ :

$$P_X = \sum_{i=1}^k w_i P_{\pi_i}.$$
(1)

Such clustering based decompositions have been considered from the information theory point of view (see, [29], [31]). Every item in the set X belongs, with certain probability, to each of the clusters so that, a clustering solution is imaged as a set of associating probability distributions. This association is termed "fuzzy membership in a cluster". A hard clustering situation appears, as a partial case, when each point is assigned to one of the clusters with the probability 1.

From the information standpoint, clustering is the basic approach for the data compression by throwing away the unrelated information. A loss compression can be made by means of assigning the data items to clusters so that the mutual information

$$I(\Pi_k, X) = \sum_{x,j} P(\pi_j \mid x) P(x) \log_2 \left( \frac{P(\pi_j \mid x)}{P(\pi_j)} \right)$$

is minimized. A difference between the relevant and irrelevant information can be provided by a distortion function  $(\cos t) E_j(x)$  evaluating the energy of assigning a point  $x \in X$  to the cluster  $\pi_j$ . The minimization

of  $I(\Pi_k, X)$  is constrained by the expected distortion

$$\overline{d}(X,\Pi_k) = \sum_{x,j} P(\pi_j \mid x) P(x) E_j(x).$$

The formal solution of this task is provided by the Boltzmann-Gibbs distributions

$$P(\pi_j \mid x) = \frac{P(\pi_j)}{Z(x,T)} \exp\left(-\frac{E_j(x)}{T}\right),$$

where

$$Z(x,T) = \sum_{j} P(\pi_{j}) \exp\left(-\frac{E_{j}(x)}{T}\right)$$

is the partition function,

 $\beta = \frac{1}{T}$  is the Lagrange multiplier and *T* is interpreted in a physical analogy for the system temperature. For  $T = \infty$ , each point is equally associated with all clusters, and T = 0 leads to a hard clustering situation.

Let us now suppose that hard clusters are defined by the vectors  $Y = \{y_i\}, j=1,...,k$  such that

$$E_{j}(x) = (x - y_{j}) \sum_{j}^{-1} (x - y_{j}),$$
(2)

where  $\boldsymbol{\Sigma}_{j}$  are the clusters covariance matrices. The marginal distribution of  $\textbf{\textit{Y}}$  becomes

$$P(Y) = e^{-\frac{F}{T}} / \sum_{Y} e^{-\frac{F}{T}},$$

where

$$F = -T\sum_{x} \ln\left(\sum_{j} \exp\left(-\frac{(x-y_j)\sum_{j}^{-1}(x-y_j)}{T}\right)\right)$$

is the so-called "the free energy" of the partition. A model closely related to one arises in clustering problems as a variant of the underlying distribution (1):

$$f(x) = \sum_{j=1}^{k} p_j G(x \mid y_j, \Sigma_j),$$
(3)

where  $G(x | y, \Sigma)$  is the Gaussian density with the mean y and the covariance matrix  $\Sigma$ . The maximumlikelihood estimation of the model parameters is provided by means of the well known *EM* algorithm. The standard k-means algorithm is its partial case (see, for example [12], [5]), when  $\Sigma_j = \sigma^2 I$ , j = 1, ..., k,  $p_1 = p_1 = ... = p_k$ , where I is the identity matrix and  $\sigma^2$  is the unknown common clusters dispersion. The essential difference between the approaches is that, in the free energy optimisation approach, no prior information except of (2) is assumed. Here, the cluster distributions are directly obtained from the suitable Boltzmann and Gibbs distributions.

## 3. The Approach Description

In the framework of our approach, we assume that the stable partition corresponds to the decomposition (3). For two i.i.d. disjoint samples  $S_1$  and  $S_2$ , independently drawn from  $X \subset \Re^d$ , and a partition  $\prod_k(X)$ , we could consider the sample occurrences in clusters:

$$S_{i,j} = S_i \cap \pi_j, i = 1, 2, j = 1, ..., k.$$

A distance between the sets  $S_{1,j}$  and  $S_{2,j}$  can be measured by the Hotelling's two-sample *T*-square statistic:

$$\hat{t}^{2}(S_{1,j}, S_{2,j}) = \frac{|S_{1,j}| |S_{2,j}|}{|S_{1,j}| + |S_{2,j}|} (\overline{S}_{1,j} - \overline{S}_{2,j}) W^{-1} (\overline{S}_{1,j} - \overline{S}_{2,j}),$$

where  $\bar{S}_{1,j}$  and  $\bar{S}_{2,j}$  are the sets means and *W* is the variance-covariance matrix. Under our presumption, actually presenting the null hypothesis according to the inner cluster distributions, the sets  $S_{1,j}$  and  $S_{2,j}$  for each j = 1,...,k could have been, consistently, drawn from the Gaussian distribution  $G(x | y_j, \Sigma_j)$ . In this case the variable  $\hat{t}^2(S_{1,j}, S_{2,j})$  has the Hotelling's *T*-square distribution  $T^2(d, |S_{1,j}| + |S_{2,j}| - 2)$  with parameters *d* and  $|S_{1,j}| + |S_{2,j}| - 2$  if  $|S_{1,j}| + |S_{2,j}| - 2 > 0$  (see, for example, [25]). In this context and under the null hypothesis, a cluster merit can be represented as the *p*-value calculated for the observed  $\hat{t}^2(S_{1,j}, S_{2,j})$ :

$$\delta_{j} = 1 - T^{2}CDF(d, \left|S_{1,j}\right| + \left|S_{2,j}\right| - 2)(\hat{t}^{2}(S_{1,j}, S_{2,j})),$$
(4)

where  $T^2CDF$  denotes the Cumulative Distribution Function of the corresponding Hotelling's *T*-square distribution. Consequently, the partition quality is given by:

$$\delta(\Pi_k) = \min_{1 \le j \le k} \delta_j.$$
<sup>(5)</sup>

I.e. the partition is characterized by the cluster which mostly differs from a Gaussian one.

Two main problems arise in the implementation of the proposed methodology. The first one is that the samples occurrences in the clusters can not be directly obtained since the required data stable partition is unknown. Moreover, according to the definition, samples involved in the calculation of the Hotelling's two-sample *T*-square statistic have to be independently drawn. The second one is caused by the fact that the corresponding *p*-values are inherently small, thus the null hypothesis can be rejected even when the true number of clusters is tested.

We handle the first problem by the procedure detailed in [1] and [35], which is proposed to simulate independent cluster occurrences in the clusters. A clustering algorithm used for this purpose is an important ingredient of the approach, because partitions constructed by means of the algorithm are intended more or less to reflect the inner hidden data structure, corresponding to the underlying mixture distributions described by (3). Cluster solutions offered by different algorithms can be essentially different.

Here, we iterate the clustering process as follows: Let *S* be a subset of *X*. A clustering algorithm  $\Delta_k$  is a function which maps *S* on  $C_k$ . Obviously, such a function supplies a partition of *S*. For a given pair of samples  $S_1$  and  $S_2$ , three partitions are introduced.

$$\begin{split} &\prod_{k,1} = \Delta_k(S_1), \\ &\prod_{k,2} = \Delta_k(S_2), \\ &\prod_k = \Delta_k(S_1 \cup S_2). \end{split}$$

Each item is located in the partition  $\Pi_k$  and in one of the other partitions  $\Pi_{k,1}$  or  $\Pi_{k,2}$ . On the other hand, a cluster can be differently marked in these partitions due to the fact that the clusters labels have no meaning. We overcome this difficulty by locating the permutations  $\psi_i$ , *i*=1, 2 of the set  $C_k$  which minimizes the misclassified quantities

$$\psi_i^* = \arg\min_{\psi} \sum_{x \in S_i} I(\psi(\alpha(\Pi_k)(x)) \neq \alpha(\Pi_{k,i})(x)), i = 1, 2,$$

where  $I(\bullet)$  is the indicator function. The Hungarian method [22] solves this problem by  $O(k^3)$  complexity. After changing the cluster labels of the partitions  $\Pi_{k,i}$ , i = 1, 2 consistent with  $\psi_i^*$ , i = 1, 2 the sets

$$S_{j,i} = \{x \in S_i \mid \alpha(\Pi_{k,i}) = j\}, i = 1, 2, j = 1, ..., k\}$$
(6)

can be considered as independent samples occurrences in the clusters.

As for the second problem, we propose to overcome this obstacle by making an inference on the number of clusters by relying on a big amount of data. Thus, a decision can be made by constricting an appropriate empirical distribution of  $\delta(\Pi_k)$  for several different numbers of clusters. It is naturally anticipated that the distribution having the shortest left tail belongs to the true number of clusters. A meta-algorithm, which implements the offered method, can be expressed as follows:

#### Algorithm

- 1. Repeat for each tested number of clusters from  $k_{min}$  to  $k_{max}$ ;
- 2. Repeat the following steps for pairs of samples randomly drawn without replacement;
- 3. Simulate the samples occurrences in the clusters;
- 4. Calculate, according to (4), the *p*-values of the distances between the occurrences of different samples within the clusters;
- 5. Calculate the partition merit, according to (5);
- 6. The estimate of the true number of clusters is the  $k^*$ , for which the distribution is the least concentrated at the origin.

# **4. Experimental Results**

This section presents the experimental results obtained from an implementation of the described methodology. First of all, let us present the used notations:

- $k_{min}$  the minimal tested number of clusters (the default value is 2);
- $k_{max}$  the maximal tested number of clusters (the default value is 8);
- *NS* the number of drawn sample pairs;
- *M* the size of the drawn samples;
- $\Delta_k$  the used clustering algorithm (the default is the standard *k*-means algorithm);

The distributions concentrations, at the origin, are characterized by two attributes:

- $P_{50}$  the sample median;
- $N_g$  the frequency of the lowest subgroup obtained by dividing the values range into g equal range subgroups (the default value of g is 30).

The results, obtained for several synthetic and real datasets, are exposed by means of comparing with the "known" source number of clusters. Sequentially, in order to exhibit the approach stability, the procedure is repeated 10 times for each datum, and the results are represented via error-bar plots of  $P_{50}$  and  $N_g$ , calculated as function of the tested number of clusters. The bar sizes equal to two standard deviations found within the trials.

## 4.1. Synthetic data

The first three simulated datasets, each having a size of 4000, are drawn from mixtures of two-dimensional Gaussian distributions, having, correspondingly, 4, 5 and 6 components and owning the same standard deviation  $\sigma = 0.35$ . The components means are located on the unit circle with equal angular neighbouring distance from each other. The datasets components overlap. The results obtained for the two first datasets for the parameters NS = 100 and M = 300 are presented on Figures 1–2.





Figure 2. Error-bar plots of  $P_{50}$  and  $N_{30}$  for the simulated five components Gaussian dataset

7

8

5

0.75

0.65 0.8 0.55

3

5

я

2

-1 L 1

2

3



Figure 3. Error-bar plots of  $P_{50}$  and  $N_{30}$  for the simulated six components Gaussian dataset with  $\sigma = 0.3$  and M = 400

Nevertheless, increasing the samples size to M = 400 for  $\sigma = 0.3$  yields the determination of the correct number of clusters (see Figure 3).

## 4.2. Real-world data

The real dataset is chosen from the text collection *http://ftp.cs.cornell.edu/pub/smart/*. This set includes three sub-collections, consisting of

- 1033 medical;
- 1460 information science;
- 1400 aerodynamics abstracts.

We select the 600 "best" terms, following the common "bag of words" method and used the data representation by means of two leading principal components. The results presented on Figure 4 show that the number of clusters is properly determined.



Figure 4. Error-bar plots of  $P_{50}$  and  $N_{30}$  for the three text collections dataset with NS = 100 and M = 100

Analogous results are observed in the situation where only 300 "best" terms are used. This dataset can be considered as a "noised" version of the previous one.

# Conclusions

In the current paper we propose a methodology for detecting the number of clusters in a data set. In the framework of our approach we assume that a stable partition corresponds to a small distance, within clusters, between pairs of i.i.d. disjoint samples. Here, the Hotelling's two-sample *T*-square statistic is selected as a distance between the pairs of sets. It is well known that the procedures for detecting the true number of clusters are noisy. Thus, conclusions have to be drawn based on a reasonable amount of information. To overcome this difficulty we construct empirical distributions of the distance between partitions, for all possible numbers of clusters, and the selected "true" number of clusters, is the one for which the distribution is the least concentrated at the origin.

To exhibit the performance of the proposed methodology we provide several numerical simulations. It is realized that it performs quite well even under an adverse situation.

# References

- 1. Barzily, Z., Volkovich, Z., Akteke-Ozturk, B., Weber, G.-W. On a minimal spanning tree approach in the cluster validation problem, *Informatica*, Vol. 2, 2009.
- Ben-Hur, A., Elisseeff, A., Guyon, I. A stability based method for discovering structure in clustered data. In: *Proceedings of Pacific Symposium on Biocomputing*, 2002, pp. 6–17.
- 3. Ben-Hur, A. and I. Guyon. Methods in Molecular Biology / M. J. Brownstein and A. Khodursky (Ed.). Humana Press, 2003, pp. 159–182.
- 4. Calinski, R. and J. Harabasz. A dendrite method for cluster analysis, *Communications in Statistics*, Vol. 3, 1974, pp. 1–27.
- 5. Celeux G., Govaert, G. A classification EM algorithm and two stochastic versions, *Computational Statistics and Data Analysis*, Vol. 14, 1992, pp. 315–332.
- 6. Chakravarthy, S. V., Ghosh, J. Scale-Based Clustering Using the Radial Basis Function Network, *IEEE Transactions on Neural Networks*, Vol. 7(5), 1996, pp. 1250–1261.
- 7. Cheng, R. and G. W. Milligan. Measuring the influence of individual data points in a cluster analysis, *Journal of Classification*, Vol. 13, 1996, pp. 315–335.
### **Applied Statistics and Operation Research**

- 8. Cuevas, A., Febrero, M., Fraiman, R. Estimating the Number of Clusters, *The Canadian Journal of Statistics*, Vol. 28 (2), 2000, pp. 367–382.
- 9. Cuevas, A., Febrero, M. and R. Fraiman. Cluster Analysis: A Further Approach Based on Density Estimation, *Computational Statistics and Data Analysis*, Vol. 28, 2001, pp. 441–459.
- 10. Dudoit S., Fridlyand, J. A prediction-based resampling method for estimating the number of clusters in a dataset, *Genome Biology*, Vol. 3(7), 2002, pp. 1–21.
- 11. Dunn, J. C. Well Separated Clusters and Optimal Fuzzy Partitions, *Journal Cybern.*, Vol. 4, 1974, pp. 95–104.
- 12. Fraley, C. and A. E. Raftery. How Many Clusters? Which Clustering Method? Answers Via Model-Based Cluster Analysis, *The Computer Journal*, Vol. 41(8), 1998, pp. 578–588.
- 13. Gordon, A. D. Identifying genuine clusters in a classification, *Computational Statistics and Data Analysis*, Vol. 18, 1994, pp. 561–581.
- 14. Gordon, A. D. Classification. Boca Raton, FL: Chapman and Hall, CRC, 1999.
- 15. Hamerly, G. and C. Elkan. Learning the *k* in *k*-means. In: *Proceedings of the Seventeenth Annual Conference on Neural Information Processing Systems (NIPS), December, 2003*, pp. 281–288.
- 16. Hartigan, J. A. Clustering Algorithms. Wiley, 1975.
- 17. Hartigan, J. A. Consistency of Single Linkage for High-Density Clusters, *Journal of the American Statistical Association*, Vol. 76, 1981, pp. 388–394.
- 18. Hartigan, J. A. Statistical theory in clustering, J. Classification, Vol. 2, 1985, pp. 63-76.
- 19. Hubert, L. and J. Schultz. Quadratic assignment as a general data-analysis strategy, *British J. Math. Statist. Psychology*, Vol. 29, 1976, pp. 190–241.
- 20. Jain, A., Dubes, R. Algorithms for Clustering Data. New Jersey: Englewood Cliffs, Prentice-Hall, 1988.
- 21. Krzanowski W., Lai, Y. A criterion for determining the number of groups in a dataset using sum of squares clustering, *Biometrics*, Vol. 44, 1985, pp. 23–34.
- 22. Kuhn, K. The Hungarian method for the assignment problem, *Naval Research Logistics Quarterly*, Vol. 2, 1955, pp. 83–97.
- 23. Lange, T., Roth, V., Braun, M., Buhmann, J. M. Stability-based validation of clustering solutions, *Neural Computation*, Vol. 15(6), 2004, pp. 1299–1323.
- 24. Levine, E., Domany, E. Resampling Method for Unsupervised Estimation of Cluster Validity, *Neural Computation*, Vol. 13, 2001, pp. 2573–2593.
- 25. Mardia, K. V., Kent, J. T. and J. M. Bibby. Multivariate Analysis. Academic Press, 1979.
- 26. Milligan, G. and M. Cooper. An examination of procedures for determining the number of clusters in a data set, *Psychometrika*, Vol. 50, 1985, pp. 159–179.
- 27. Mufti, G., Bertrand, P., El Moubarki, L. Determining the number of groups from measures of cluster validity. In: *Proceedigns of ASMDA*, 2005, pp. 404–414.
- Pelleg, D., Moore, A. X-means: Extending K-means with efficient estimation of the number of clusters. In: *Proceedings of the 17th International Conf. on Machine Learning*. San Francisco, CA: Morgan Kaufmann, 2000, pp. 727–734.
- 29. Rose, K., Gurewitz, E. and G. Fox. Statistical mechanics and phase transitions in clustering, *Physical Review Letters*, Vol. 65(8), 1990, pp. 945–948.
- 30. Roth, V., Lange, T., Braun, M. and J. Buhmann. A resampling approach to cluster validation. In: *Proc. Intl. Conf. on Computational Statistics*, 2002, pp. 123–128.
- 31. Still, S. and W. Bialek. How Many Clusters? An Information-Theoretic Perspective, *Neural Computation*, Vol. 16(12), 2004, pp. 2483–2506.
- 32. Stuetzle, W. Estimating the Cluster Tree of a Density by Analyzing the Minimal Spanning Tree of a Sample, *J. Classification*, Vol. 20(5), 2003, pp. 25–47.
- 33. Sugar, C. and G. James. Finding the Number of Clusters in a Data Set: An Information Theoretic Approach, *Journal of the American Statistical Association*, Vol. 98, 2003, pp. 750–763.
- 34. Tibshirani, R., Walther, G. and T. Hastie. Estimating the number of clusters via the gap statistic, *J. Royal Statist. Soc. B*, Vol. 63(2), 2001, pp. 411–423.
- 35. Volkovich, Z., Barzily, Z. and L. Morozensky. A statistical model of cluster stability, *Pattern Recognition*, Vol. 41(7), 2008, pp. 2174–2188.
- 36. Volkovich, Z., Barzily, Z., Avros, R. and D. Toledano-Kitai. On application of the *K*-nearest neighbors approach for cluster validation. In: *Proceedings of the XIII International Conference Applied Stochastic Models and Data Analysis (ASMDA)*, 2009.
- Wishart, D. Mode Analysis: A Generalization of Nearest Neighbor, which Reduces Chaining Effects. In: *Numerical Taxonomy*. London: Academic Press, 1969, pp. 282–311.

Received on the 21st of June, 2010

Computer Modelling and New Technologies, 2010, Vol.14, No.4, 73–76 Transport and Telecommunication Institute, Lomonosova 1, LV-1019, Riga, Latvia

### PREDICTIONS ON THE FORMATION OF NERVE-MUSCLE CONNECTION

### I. Nowik

The Department of Industrial Engineering and Management, SCE-Shamoon College of Engineering Beer-Sheva, 84100, Israel E-mail: iritno@sce.ac.il

During the first couple of weeks after birth neurons that innervate a muscle "compete" with each other to innervate a maximal number of muscle fibres. The result of this competition is that the less active neurons innervate more muscle fibres than the more active neurons. This is called the "size principle" and it means that the less active neurons win in more competitions. This is surprising, as looking at the *isolated* muscle fibre (i.e., one battle) a competitive advantage to the *more* active neurons is observed. Using the approach of game theory we explain, in a previous paper [1], how this may happen. In the current paper we present several predictions that follow from our model.

Keywords: game theory; size principle; computational neuroscience

### **1. Introduction**

A muscle is composed of many fibres. At birth, each muscle fibre is innervated by several neurons called motoneurons (MNs) but during the first couple of weeks after birth, a competitive process called "synapse elimination" abolishes the connections of all MNs but one, which we call "the winner at that muscle fibre". As each MN innervates initially many muscle fibres, it engages in many competitions simultaneously, winning at some and losing at others. When the competition period ends each muscle fibre is innervated by only one MN, but each MN innervates a group of muscle fibres, called a "muscle unit". When an electrical stimulus arrives at the MN's cell body, if it is higher than the MN's activation threshold, then it activates the MN, which in turn, activates all the fibres in its muscle unit. At the end of synapse elimination, the less active MNs (i.e., MNs with higher activation thresholds) have larger muscle units than the more active MNs. This is called "The Size Principle". The size principle is well established empirical fact. It is found in the motor system of almost all vertebrates, and is thought of as one of the most fundamental principles in the organization of motor-unit behavior, therefore it is important to understand how it evolves. In viewing the period of synapse elimination as a game in which MNs are competing to innervate a maximal number of muscle fibres, the translation of the size principle is that the less active MNs win in more competitions than the more active MNs. This means that being less active is "advantageous" in this process. But surprisingly, the majority of experiments that have manipulated the activity of MNs during synapse elimination seem to point to the opposite conclusion, that it is the more active MNs that are advantageous in this process. In particular, all the experiments that were done at the isolated muscle fibre suggest some competitive advantage to the more active MNs. Thus, the less active MNs are disadvantage in the single battles (i.e., single muscle fibre) but manage, somehow, to win the war (i.e., win in more muscle fibres). In this work we explain why the less active MNs win in more competitions and we resolve the paradox of seemingly contradicting experimental data. The results of this work are proved mathematically and are presented in a separate paper [1]. In order to illustrate these results we additionally simulated the game using *MatLab*. In the present paper, we concentrate on several predictions that follow from our model.

### 2. The Model

The players are the MNs innervating one muscle. As we are interested in the connection between the activity level of a MN and the size of it's muscle units, we define the *strategy* of each MN as its activity level,<sup>\*</sup> and define the *payoff* as the final size of its muscle unit. We assume that the activity levels

As in evolutionary games, it is not assumed that players choose their strategies, but rather they are determined by the genes. Unlike evolutionary games, the activity levels do not change along the game, thus from an evolutionary game theory stand point this is a one stage game.

### **Applied Statictics and Operation Research**

of the MNs are independent identically distributed random variables.<sup>†</sup> We divide the MNs into two equal sized groups according to their activity levels: M-group contains the 50% more active MNs and L-group contains the 50% less active MNs. This division does not imply that MNs in the same group cooperate in any manner (e.g., share resources). It is only done for mathematical convenience.

We would like to define the activity level of a *muscle fibre*, namely the frequency of its activation. As the MNs connecting to the muscle fibre are the ones to activate it, we define the level of activity of a muscle fibre as the sum of activity levels of the MNs connecting to it.

There are two rules in the game, based on experimental results:

- 1. The competitions at the muscle fibres end at different times, according to a decreasing level of activity of the muscle fibres, namely the first competition to end, is the one occurring at the most active muscle fibre and the last competition to end occurs at the least active muscle fibre.
- 2. When a MN wins a muscle fibre this reduces its future winning probabilities at other muscle fibres.

The main idea of our model and the key factor in understanding why the less active MNs win in more competitions is realizing that the *time* of winning at a muscle fibre plays a critical role in the process. When a MN wins at a muscle fibre it must devote resources for maintaining this connection, and thus it has less available resource for competing at other muscle fibres. In such circumstances it is advantageous to win in *later* stages of the process rather than in earlier ones, because winning at a late stage will affect only the fewer competitions that are not yet resolved (i.e., the MN will lack resources and compete poorly only in the few remaining competitions), whereas winning at earlier stages will negatively affect more competitions and cause the MN to lose in more competitions. This is exactly what the less active MNs do; they mainly take part in competitions that end at later stages of the game and thus, these are the competitions in which they win. The reason for our claim that the less active MNs take part in competitions that end relatively late is as follows. According to rule 1 the time in which a competition ends is determined according to the activity level of the muscle fibre but this, in turn, is determined according to the activity level of the MNs connecting to it. Therefore a muscle fibre that is mainly connected by more active MNs will be more active and thus its competition will end at early stage of the process. Similarly, a muscle fibre that is mainly connected by less active MNs will be less active and thus its competition will end at later stages of the process. Summarizing this argument - the more active MNs tend to win at earlier stages of the process and the less active MNs tend to win in later stages of the process. This difference in the times of winning works in favour of the less active MNs, because, as explained, it is advantageous to win later rather than earlier. As a result of this advantage, the less active MNs win, in total, more competitions than the more active MNs. Thus even though, the less active MNs may have some disadvantage at the single battle (single muscle fibre), because they "invest" and win more in *later* competitions, then in total, they win in more competitions and thus 'win the war' yielding the size principle.

We now present several predictions that follow from our model. We use *MatLab* to illustrate some of these predictions.

### **3. Predictions**

The predictions presented here are of two types. The first, are predictions that relate to known experimental results. In such case we show that our model is consistent with these results. The second type, are *new* predictions that follow from our model.

# **3.1.** The MN Winning at a Muscle Fibre is Expected to be Less and Less Active as Synapse Elimination Proceeds

This prediction follows directly from the rule 1 listed above.

### 3.2. The Role of Neuronal Indentity in Synapse Elimination

In a paper by Kasthuri and Lichtman [2], all muscle fibres that were co-innervated by the same two MNs were examined on a late stage of synapse elimination (between the seventh to ninth days after birth). At this late stage of the game, these two MNs were usually the last two competitors left at these

<sup>&</sup>lt;sup>†</sup> The distribution density is assumed to be symmetrical, e.g., uniform or normal distribution.

### **Applied Statictics and Operation Research**

muscle fibres. The competitions all ended with the same winner. This winner was the one that had a smaller muscle unit at that time. Our model predicts this experimental result according to rule 2.

#### **3.3. Selective Stimulation**

Ridge and Betz [3] selectively stimulated the activity of some of the MNs innervating a muscle. This resulted in larger muscle units for the stimulated group. According to our model, selective stimulation is expected to have opposite effects; raising the winning probabilities of the stimulated MNs during the competition period, but also specifically bringing forward competitions at muscle fibres that are innervated by stimulated axons (thus reducing their actual winning probabilities).

To introduce stimulation in our model, the MNs were divided to "stimulated" and "un-stimulated" groups (instead of "more active" and "less active" MNs). Figure 1a shows a simulation of our model that follows the stimulation procedure of Ridge and Betz [3]. As seen, our model yields the same result as in the experimental results. In addition, we predict that executing the same procedure *earlier* will prove to be less successful for the stimulated group (Figure 1b) (see [1] for details regarding the simulation procedure).



relative stage of the game

Figure 1. The effect of selective stimulation according to our model.

a) As in Ridge & Betz[3], 5 consecutive short stimulations were applied, each for a fraction of 0.015 stages (corresponding to 4 hrs of stimulation per day, for 5 days). Depicted, is the difference in the number of winnings between the stimulated and the un-stimulated groups. The stimulated group (dotted line) won in significantly more competitions than the control (smooth line) P < 0.03, one-tailed t-test.</li>
 b) Executing the same procedure earlier is less successful for the stimulated group

### **Applied Statictics and Operation Research**

### 3.4. The Size of the Muscle is not a Factor in the expression of the Size Principle

We measure the expression of the size principle by the difference in the number of wins (between the M-group and L-group), divided by the number of fibres in the muscle. It follows from the mathematical analysis of our model (see [1]) that a large muscle (i.e., a muscle with many fibres) expresses the size principle to the same extent as a small muscle. This is also shown in the simulation presented on Figure 2.



*Figure 2.* According to our model the expression of the size principle does not depend on the size of the muscle. Depicted is the difference W between the number of wins of the two groups, divided by the number of fibres N in the muscle. As seen, All along the game, there is no significant difference in the expression of the size principle for different values of N

### Conclusions

In this paper we presented several predictions that follow from our model for the competition between MNs innervating one muscle. Although the assumptions of the game are based on experiments, they are easily generalized for application to a much wider scope of situations, for example economical systems. When facing a multi-stage 'game', in which resources are limited and are needed for maintenance of previous wins, winning later is advantageous, thus the strategic implication is that if one needs to allocate his resources in advance, one should invest more in later competitions rather than in earlier ones.

### References

- 1. Nowik, I. The game MNs play, *Games and Economic Behavior*, Vol. 66, 2009, pp. 426–461.
- 2. Kasthuri, N., Lichtman, J. W. The role of neuronal identity in synaptic competition, *Nature*, Vol. 424, 2003, pp. 426–30.
- 3. Ridge, R. M. and W. J. Betz. The effect of selective, chronic stimulation on motor unit size in developing rat muscle, *J. Neurosci*, Vol. 4, 1984, pp. 2614–20.

Received on the 21st of June, 2010

Computer Modelling & New Technologies, 2010, Volume 14, No. 4 Transport and Telecommunication Institute, Lomonosova 1, Riga, LV-1019, Latvia

# Authors' index

Avros R.	65
Barzily Z.	65
Ben-Yair A.	14
Bischoff W.	7
Golenko-Ginzburg D.	50
Greenberg D.	14
Greenglaz L.	19
Gurevich G.	6,31
Guseynov S.	19
Huber C.	40
Iliev V.A.	57
Kopytov E.	19
Laslo Z.	6,50
Nowik I.	73
Petkov G.I.	57
Puzinkevich E.	19
Toledano-Kitai D.	65
Vexler A.	31
Volkovich Z.	65
Vonta F.	40

### Computer Modelling & New Technologies, 2010, Volume 14, No. 4 \*\*\* Personalia



### Yuri N. Shunin (born in Riga, March 6, 1951)

- Vice-Rector on Academic Issues (Information Systems Management Institute), professor, Dr.Sc.Habil., Member of International Academy of Refrigeration
- Director of Professional Study Programme Information Systems (Information Systems Management Institute)
- Director of Master Study Programme Computer systems (Information Systems Management Institute)
- University study: Moscow physical and technical institute (1968–1974).
- Ph.D. (physics & mathematics) on solid state physics (1982, Physics Institute of Latvian Academy of Sciences), Dr.Sc.Habil (physics & mathematics) on solid state physics (1992, Ioffe Physical Institute of Russian Academy of Sciences)
- Publications: 400 publications, 1 patent
- Scientific activities: solid state physics, physics of disordered condensed media, amorphous semiconductors and glassy metals, semiconductor technologies, heavy ion induced excitations in solids, mathematical and computer modelling, system analysis



- Vice-rector for Research and Development Affairs of Transport and Telecommunication Institute, Professor, Director of Telematics and Logistics Institute
- PhD in Aviation (1981, Moscow Institute of Civil Aviation Engineering) Dr.Sc.Habil. in Aviation (1992, Riga Aviation University), Member of the International Telecommunication Academy, Member of IEEE, Corresponding Member of Latvian Academy of Sciences (1998)
- Publications: 420 scientific papers and 67 patents
- **Research activities:** information technology applications, operations research, electronics and telecommunication, analysis and modelling of complex systems, transport telematics and logistics



#### Filia Vonta

- Assistant Professor, Department of Mathematics, National Technical University of Athens, Greece
- Education-employment: she received a MA in 1988 and a PhD degree in 1992 from the University of Maryland, USA, in Mathematical Statistics. From 1993 to 2009 she was employed at the Department of Mathematics and Statistics of the University of Cyprus, Cyprus
- Scientific interests: semiparametric Statistics, Survival analysis, Analysis of Medical Data
- **Publications:** more than 30 in refereed journals and books



### Doron Greenberg (born in Israel, 1955)

- Senior Lecturer and Head of the Financial Branch, the Department of Economics and Business Administration and the Department of Industrial Engineering, Ariel University Center of Samaria, Ariel, Israel
- University study: Leon Recanati Graduate School of Business Administration, Tel-Aviv University, Israel (1982–1985)
- University study: The Technion, Israel Institute of Technology, Haifa, Israel (1978–1980)
- Ph.D. (Economics) on applying option theory to investments in R&D (1992, the University of Houston, USA)
- **Publications:** about 30 refereed articles and refereed letters in scientific journals
- Scientific activities: economic capital risk management, promoting ethics in organizations, production planning and control, planning and controlling network projects, industrial scheduling, managing reliability and safety

### Computer Modelling & New Technologies, 2010, Volume 14, No. 4 \*\*\* Personalia



### Dimitri Golenko-Ginzburg (born in Moscow, November 24, 1932)

- Professor, Industrial Engineering and Management Department, Ariel University Center of Samaria, Ariel, Israel
- Professor, Industrial Engineering and Management Department, Ben-Gurion University of the Negev, Beer-Sheva, Israel (1988-2004)
- Full Professor (tenured position), Institute of National Economics, Uzbekistan Ministry of Higher Education, Tashkent (1977–1979)
- Full Professor (tenured position), Moscow Economico-Statistical Institute, USSR Ministry of Higher Education, Moscow (1967–1977)
- University study: Moscow State University, Department of Mathematics (1954–1958)
- University study: Moscow Institute of National Economics, Department of Economics (1950–1954)
- Ph.D. (Applied Mathematics) on simulating probability processes on computers (1962, Moscow Physico-Technical Institute)
- Publications: 15 books, about 500 refereed articles and refereed letters in scientific journals
- Scientific activities: production planning and control, planning and controlling network projects, industrial scheduling, managing reliability and safety

#### Avner Ben-Yair (born in Moscow, May 19, 1961)

- Lecturer, Industrial Engineering and Management Department, SCE Sami Shamoon College of Engineering, Beer-Sheva, Israel (2002–2010)
- Health and Safety-at-Work Engineer, Baran Engineering Co., Israel (1996–1999)
- Health and Safety-at-Work Engineer, AVX Israel, Jerusalem (1986-1996)
- University study: Ben-Gurion University of the Negev, Beer-Sheva, Israel (1999–2001)
- University study: Moscow Polygraphic Institute, Department of Mechanical Engineering (1979–1985)
- Ph.D. (Industrial Engineering and Management) on harmonization models in strategic management and safety engineering (2004, Ben-Gurion University of the Negev)
- Publications: 90 refereed articles and refereed letters in scientific journals
- Scientific activities: system performance and effectiveness, reliability and failure analysis, fault tree analysis, trade-off optimisation models for organization systems, risk analysis and contingency planning, maintainability and hazard analysis techniques, economic aspects of safety, production planning, scheduling and control, strategic management, network models structure and project scheduling, cost optimization and PERT-COST models



•

## **CUMULATIVE INDEX**

## COMPUTER MODELLING and NEW TECHNOLOGIES, volume 14, No. 4, 2010 (Abstracts)

**W. Bischoff.** Residual Partial Sums Techniques for Fixed Designs to Find Change-Points in Linear Regression, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 7–13.

We investigate a data set describing the quality of a production process. By the information of these data it has to be decided whether the quality is constant or whether the quality changes. Our null hypothesis is that the quality is constant that is a linear regression. In practice it is popular to investigate the partial sums of the least squares residuals to look for changes in linear regression. The partial sums of the least squares residuals can be embedded into the class of continuous functions. By this procedure we obtain a stochastic process with continuous paths. It is called residual partial sum process. If the number of observations is large enough a projection of the Brownian motion can be considered as approximation (with respect to weak convergence) of the residual partial sum process. This projection of the Brownian motion can be used to establish non-parametric tests of Cramér-von Mises and Kolmogorov-Smirnov type to test for changes in linear regression. We use this procedure to test the data for constant quality.

**Keywords:** residual partial sum limit processes, linear regression models, fixed designs, Brownian motion, projections of Brownian motion, reproducing kernel Hilbert space, change-point problem

**D. Greenberg, A. Ben-Yair.** Beta-Distribution Models in Stochastic Project Management, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 14–18.

A research is undertaken to justify the use of beta-distribution p.d.f. for man-machine type activities under random disturbances. The case of using one processor, i.e., a single resource unit, is examined. It can be proven theoretically that under certain realistic assumptions the random activity – time distribution satisfies the beta p.d.f. Changing more or less the implemented assumptions, we may alter to a certain extent the structure of the p.d.f. At the same time, its essential features (e.g. asymmetry, unimodality, etc.) remain unchanged. The outlined above research can be applied to semi-automated activities, where the presence of man-machine influence under random disturbances is, indeed, very essential. Those activities are likely to be considered in organization systems (e.g. in project management), but not in fully automated plants.

**Keywords:** random activity duration, time – activity beta-distribution, operating by means of a single processor, convergence to a beta-distribution "family"

**E. Kopytov, S. Guseynov, E. Puzinkevich, L. Greenglaz.** Continuous Models of Current Stock of Divisible Productions, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 19–30.

In the given paper we investigate the problem of constructing continuous and unsteady mathematical models to determine the volumes of current stock of divisible productions using apparatus and equations of mathematical physics. It is assumed that time of production distribution and replenishment is continuous. The constructed models are stochastic, and have different levels of complexity, adequacy and application potentials. The simple models are constructed using the theory of ordinary differential equations, for construction of more complex models the theory of partial differential equations is applied. Furthermore for some of proposed models we have found an analytical solution in the closed form, and for some of proposed models the discretization is carried out using stable difference schemes.

**Keywords**: *inventory control model, current stock, divisible production, equations of mathematical physics* 

**G. Gurevich, A. Vexler.** Statistical Inference Using Entropy Based Empirical Likelihood Statistics, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 31–39.

In this article, we show that well known entropy-based tests are a product of empirical likelihood ratio. This approach yields stable definitions of entropy-based statistics for goodness-of fit tests and provides a simple development of two-sample tests based on samples entropy that have not been presented in the literature. We introduce the distribution-free density-based likelihood techniques, applied to test for goodness-of-fit. In addition, we propose and examine nonparametric two-sample likelihood ratio tests for the case-control study based on samples entropy. The Monte Carlo simulation study indicates that the proposed tests compare favourably with the standard procedures, for a wide range of null/alternative distributions.

**Keywords:** *empirical likelihood, entropy, goodness-of-fit tests, two-sample nonparametric tests, case-control study* 

**F. Vonta**, **C. Huber.** On the Estimation of Structural Parameters in Frailty Models for Interval Censored and Truncated Data, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 40–49.

We consider survival data that are both interval censored and interval truncated. We assume a semiparametric frailty or transformation model for the survival function and consider censoring and truncation distributions as in Huber, Solev and Vonta [6], [7]. We propose the use of modified profile likelihood estimators for the structural parameter of the model as in Slud and Vonta [11]. For fixed values of the structural parameter, we derive the least favourable parametrization of the nuisance infinite-dimensional parameter, on which the definition of the modified profile likelihood estimator is relied upon. We discuss the semiparametric efficiency of the modified profile likelihood estimator of the finite-dimensional regression parameter in the presence of the infinite-dimensional nuisance parameter, that is, the baseline cumulative hazard function.

**Keywords:** *frailty models, least favourable model, interval censored and truncated data, semiparametric estimation* 

**D. Golenko-Ginzburg, Z. Laslo.** Upon Controlling Several Building Projects in a Two-Level Construction System, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 50–56.

A two-level construction system is considered to be composed of several different building projects  $U_i$ ,  $1 \le i \le n$ , at the lower level and a control device at the upper one. The upper system's level is required to produce a given target amount V by a given due date D subject to a chance constraint, i.e. the least permissible probability P of meeting the target on time is pregiven. Each building project  $U_i$ 

has several possible speeds  $v_{i1}$ ,  $v_{i2}$ , ...,  $v_{im}$ , which are subject to random disturbances. The project's output can be measured only at preset inspection (control) points. The target amount is gauged by a single measure, e.g. in square meters, and may be rescheduled among the projects. For each unit, the average costs per time unit for each project and the average cost of performing a single inspection at a control point to observe the actual output at that point are given.

We present a two-level on-line control model under random disturbances, which centres on minimizing the system's expenses subject to the chance constraint. The suggested two-level heuristic algorithm is based on rescheduling the overall target among the projects both at t = 0, when the system starts functioning, and at each emergency point, when it is anticipated that a certain project is unable to meet its local target on time subject to a chance constraint. At any emergency point t the remaining system's target  $V_t$  is rescheduled among the projects; thus, new local targets  $V_{it}$ ,  $1 \le i \le n$ ,  $\sum_i V_{it} = V_t$ , are determined. New local chance constraint values  $p_{it}$  are determined too. Those values enable

the system to meet its overall target at the due date subject to the pregiven chance constraint p.

**Keywords:** production speed, cost-optimisation, target amount reassignment, chance constraint, inspection point

**G. I. Petkov, V. A. Iliev.** Performance Evaluation of Teamwork Reflexive Analysis, *Computer Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 57–64.

The paper presents the development and supplementation of an advanced Human Reliability Analysis technique – Performance Evaluation of Teamwork method with approaches and mechanisms of reflexive analysis for prevention of erroneous *communication and decision-making*.

**Keywords:** reflection, reflexive analysis, typing, scenario analysis, styling, communicative barriers, Human Reliability Assessment (HRA), risk, human erroneous actions, Performance Evaluation of Teamwork (PET) method

**Z. Volkovich, Z. Barzily, D. Toledano-Kitai, R. Avros.** The Hotelling's Metric as a Cluster Stability Measure, *Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 65–72.

The estimation of the suggested number of clusters in dataset is an ill posed problem of essential relevance in cluster analysis. A group (cluster) is characterized by a relatively high similarity among its elements in addition to a relatively low similarity to elements of other groups. High stability in partitions, obtained from the same data source, is logically classified as a high consistency of the clustering process. Thus, the number of clusters that maximizes cluster stability can serve as an estimator for the "true" number of clusters. In the current paper we consider a probabilistic approach to this problem resting upon the Gaussian clusters model. We claim that sequences of clustered samples can be interpreted as Gaussian distributed i.i.d. samples drawn from the same source, if the number of clusters is chosen correctly. The samples closeness, within the clusters, can be measured by means of the *p*-values, calculated for the appropriate Hotelling's *T*-square statistic. Data outliers and clustering algorithm's shortcomings can yield small *p*-values. However, we presume that their empirical distribution is least concentrated at the origin, if the number of clusters is chosen correctly. Our procedure can roughly be described as a generation of such a distribution with a sequential application of a *concentration test*.

Keywords: clustering, Hotelling's criteria, cluster validation

**I. Nowik.** Predictions on the Formation of Nerve-Muscle Connection, *Modelling and New Technologies*, vol. 14, No 4, 2010, pp. 73–76.

During the first couple of weeks after birth neurons that innervate a muscle "compete" with each other to innervate a maximal number of muscle fibres. The result of this competition is that the less active neurons innervate more muscle fibres than the more active neurons. This is called the "size principle" and it means that the less active neurons win in more competitions. This is surprising, as looking at the *isolated* muscle fibre (i.e., one battle) a competitive advantage to the *more* active neurons is observed. Using the approach of game theory we explain, in a previous paper [1], how this may happen. In the current paper we present several predictions that follow from our model.

**Keywords:** game theory, size principle, computational neuroscience

## COMPUTER MODELLING and NEW TECHNOLOGIES, 14.sējums, Nr. 4, 2010 (Anotācijas)

V. Bišofs. Atlikušo daļējo summu tehnikas nemainīgiem projektiem, lai atrastu maiņas punktus lineārajās regresijās, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 7.–13. lpp.

Rakstā tiek izpētīta datu rinda, kas apraksta ražošanas procesa kvalitāti. Saskaņā ar šo datu sniegto informāciju, ir jānolemj, vai kvalitāte ir konstanta, vai tā mainās. Mūsu nulles hipotēze ir, ka kvalitāte ir konstanta, kas ir lineāra regresija. Praksē ir ierasts izpētīt daļējās summas no mazākās kvadrātu starpības, lai redzētu izmaiņas lineārajā regresijā. Daļējās summas no mazākās kvadrātu starpības var tikt ievietotas nepārtraukto funkciju klasē. Līdz ar šo procedūru mēs iegūstam stohastisko procesu ar nepārtrauktiem ceļiem. Tas tiek saukts par atlikušo daļējo summas procesu. Ja novērojumu skaits ir pietiekami liels, Brauna kustības projekcija var tikt uzskatīta kā atlikušā daļējās summas procesa aproksimācija (pieņemot, ka konverģence ir vāja). Šī Brauna kustības projekcija var tikt pielietota, lai noteiktu *Cramér-von Mises* neparametriskos testus un *Kolmogorova-Smirnova* tipa testus lineārās regresijas izmaiņām. Mēs lietojam šo procedūru, lai testētu datus konstantai kvalitātei.

**Atslēgvārdi:** atlikušās daļējās summas limita procesi, lineārās regresijas modeļi, nemainīgie projekti, Brauna kustība, Brauna kustības projekcija, maiņas punkta problēma

**D. Grīnbergs**, **A. Ben-Jears**. Beta sadales modeļi stohastiskajā vadības projektā, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 14.–18. lpp.

Pētījumā tiek pamatots beta-sadales *p.d.f.* lietojums cilvēka-ierīces tipa aktivitātēm pie nejaušiem traucējumiem. Tiek izskatīts viena procesora lietojuma gadījums, t.i., viena vienīga resursa vienība. Var tikt teorētiski pierādīts, ka pie noteiktiem reāliem pieņēmumiem nejauša darbība – laika sadalījums apmierina beta *p.d.f.* Vairāk vai mazāk pamainot ieviestos pieņēmumus, mēs varam zināmā mērā pārveidot *p.d.f.* struktūru. Tanī pašā laikā tā pamata īpašības (e.g. asimetrija, uni-modalitāte, etc.) paliek nemainīgas. Iepriekšminētais pētījums var tikt pielietots semi-automatizētām darbībām, kur cilvēka-ierīces ietekmes klātbūtne pie nejaušiem traucējumiem, bez šaubām, ir ļoti būtiska. Šīm darbībām jābūt izskatītām organizācijas sistēmās (e.g. projekta vadībā), bet ne pilnīgi automatizētās rūpnīcās.

**Atslēgvārdi:** *nejaušas darbības turpināšanās, laiks – darbības beta-sadalījums, darbojoties ar viena vienīga procesora līdzekļiem, konverģence beta-sadalījuma 'kopā'* 

**J. Kopitovs, S. Guseinovs, E. Puzinkevičs, L. Gringlazs.** Dalāmās produkcijas kārtējo sortimentu nepārtrauktie modeļi, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 19.–30. lpp.

Dotajā rakstā autori pēta nepārtrauktu un nestabilu matemātisko modeļu uzbūvi, lai noteiktu dalāmās produkcijas kārtējo sortimentu, lietojot matemātiskās fizikas vienādojumus un ierīces. Tiek pieņemts, ka produkcijas sadales laiks un atkārtota piepildīšana ir nepārtraukta. Uzbūvētie modeļi ir stohastiski un tiem ir sarežģītības, atbilstības un pielietojumu potenciāla dažādi līmeņi. Vienkāršie modeļi ir uzbūvēti, pielietojot vienkāršu diferenciālvienādojumu teoriju, sarežģītāku modeļu uzbūvei tiek pielietota daļēja diferenciālvienādojumu teorija. Turpmāk dažiem no piedāvātajiem modeļiem autori ir atraduši analītisku risinājumu slēgtā veidā un dažiem no piedāvātajiem modeļiem tiek izstrādāta diskretizācija, lietojot stabilas atšķirības shēmas.

Atslēgvārdi: inventāra kontroles modelis, kārtējais krājums, dalāmā produkcija, matemātiskās fizikas vienādojumi

**G. Gurevičs, A. Vekslers.** Statistiskās metodes, lietojot entropiju, pamatotu uz empīriskās varbūtības statistiku, *Computer Modelling and New Technologies,* 14.sēj., Nr.4, 2010, 31.–39. lpp.

Šajā rakstā autori parāda, ka vispārzināmie uz entropiju bāzētie testi ir empīriskās varbūtības proporcijas produkts. Šī pieeja prasa stabilas uz entropiju bāzētas statistikas definīcijas labskanības testiem un nodrošina divu paraugu testu vienkāršu izstrādi, kas pamatoti uz paraugu entropiju, kas nav līdz šim prezentēti literatūrā. Autori piedāvā sadales brīvu uz blīvumu bāzētu varbūtības tehnikas, kas pielietots labskanības testēšanai. Papildus autori piedāvā un novērtē neparametriskus divu paraugu varbūtības proporciju testus gadījuma-kontroles izpētei, kas pamatota uz paraugu entropiju. Monte Karlo simulācijas izpēte norāda, ka piedāvātie testi salīdzināmi ar standarta procedūrām plašai nulles/alternatīvas sadales rindai.

Atslēgvārdi: empīriskā varbūtība, entropija, labskanības testi, neparametriskie divu paraugu testi, gadījuma-kontroles izpēte

**F. Vonta, K. Hubers.** Par strukturālo parametru novērtēšanu trausluma modeļos intervāla cenzētajiem un nogrieztajiem datiem, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 40.–49. lpp.

Autori izskata paliekošos datus, kuri ir kā intervāla cenzētie, tā arī intervāla nogrieztie dati. Tiek pieņemts semi-parametrisks trausluma vai transformācijas modelis izdzīvošanas funkcijai un tiek izskatītas cenzētās un nogrieztās distribūcijas, kā tas ir parādīts, skat. Huber, Solev un Vonta [6], [7]. Autori piedāvā modificētā profila varbūtības novērtētājus modeļa strukturālajam parametram, kā tas ir parādīts, skat. Slud un Vonta [11].

Strukturālā parametra fiksētām vērtībām mēs iegūstam vismaz izdevīgu parametrizāciju no neierobežoti-dimensionāla parametra traucējuma, uz kā balstās modificētā profila varbūtības novērtējuma definīcija.

Atslēgvārdi: trausluma modeļi, vismazākā izdevīguma modelis, intervāla cenzētie un nogrieztie dati, semi-parametriskā novērtēšana

**D. Golenko-Ginzburgs, Z. Laslo.** Par dažu ēku projektu kontroli divlīmeņa būves sistēmā, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 50.–56. lpp.

Divlīmeņa būves sistēma — tā sastādās no vairākiem dažādiem projektiem  $U_i$ ,  $1 \le i \le n$ , zemākajā līmenī un kontroles ierīce uz augstākā. Augstākās sistēmas līmenim ir paredzēts veikt dotā uzdevuma daudzumu V ar doto maksājumu datumu D atkarībā no gadījuma ierobežojuma, i.e. tiek iepriekš dots paredzamais uzdevums laikā vismazākajā atļaujamajā varbūtībā p. Ikvienam ēkas

projektam  $U_i$  ir vairāki iespējamie ātrumi  $v_{i1}$ ,  $v_{i2}$ , ...,  $v_{im}$ , kuri ir nejaušo dislokāciju subjekti. Projekta iznākums var būt mērīts tikai pašreizējās inspekcijas (kontroles) punktos. Uzdevumu daudzums tiek kalibrēts ar vienkāršu mēru, piem., kvadrātmetros, un var tikt pārplānots starp projektiem. Katrai vienībai tiek dotas vidējās izmaksas uz laika vienību katram projektam un vidējās izmaksas, lai veiktu vienkāršu inspekciju kontroles punktā, lai novērotu faktisko iznākumu punktā.

Autori prezentē divlīmeņu tiešsaistes kontroles modeli pie nejaušiem traucējumiem, kas centrējas uz izmaksu samazināšanos atkarīgas no gadījuma ierobežojumiem. Piedāvātais divlīmeņu heiristiskais algoritms ir pamatots uz vispārējo uzdevuma pārplānošanu starp projektiem kā pie t = 0, kad sistēma uzsāk darbību, tā pie katra kritiskā stāvokļa punkta, kad tas tiek sagaidīts, ka konkrēts projekts nespēj pārvarēt tā vietējo mērķi laikā pakļautu gadījuma ierobežojumam. Katrā kritiskā stāvokļa punktā t atlikušais sistēmas mērķis V<sub>t</sub> tiek pārplānots starp projektiem; tādējādi jauni lokālie mērķi  $V_{it}$ ,  $1 \le i \le n$ ,  $\sum_i V_{it} = V_t$ , tiek noteikti. Bez tam tiek noteiktas arī jaunas lokāla gadījuma ierobežojuma vērtības  $p_{it}$ . Šīs vērtības sekmē sistēmai pārvarēt tā vispārējo uzdevumu saskaņā ar maksājumu datumu atkarībā no iepriekš dotā gadījuma ierobežojuma p.

Atslēgvārdi: ražošanas ātrums, izmaksu optimizācija, uzdevumu daudzuma reasignējums, gadījuma ierobežojums, inspekcijas punkts

**G. I. Petkovs, V. A. Iļjevs.** Komandas darba refleksīvās analīzes veikuma novērtēšana, *Computer Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 57.–64. lpp.

Rakstā autori parāda progresīvo *Human Reliability Analysis* tehniku – komandas darba veikuma novērtēšanas metode ar refleksīvās analīzes pieejām un mehānismiem, lai aizkavētu kļūdaino komunikāciju un lēmumu pieņemšanu.

**Atslēgvārdi:** refleksija, refleksīvā analīze, cilvēka kļūdainās darbības, cilvēka uzticamības izvērtēšana (Human Reliability Assessment (HRA)), komandas darba veikuma novērtēšanas metode

**Z. Volkovičs, Z. Barzili, D. Toledano-Kitai, R. Avros.** Birojviesošanas metrika kā klastera stabilitātes mērs, *Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 65.–72. lpp.

Piedāvātā klasteru skaita datu kopās novērtēšana ir nepareizi ierosināta būtiskas saistības problēma klasteru analīzē. Grupa (klasters) tiek raksturota ar relatīvi augstu līdzību starp tās elementiem, turklāt elementiem no citām grupām ir relatīvi zema līdzība. Augsta stabilitāte nodalījumos, kas iegūta no tā paša datu avota, tiek loģiski klasificēta kā klasteringa procesa augsts nepretrunīgums. Tādējādi klasteru skaits, kas palielina klasteru stabilitāti, var kalpot kā novērtētājs patiesam klasteru skaitam. Dotajā rakstā autori izskata varbūtības pieeju šai problēmai, pamatojoties uz Gausa klasteru modeli. Autori uzskata, ka klasteru paraugu sekvences var būt interpretētas kā Gausa sadalītie i.i.d. paraugi, izvilkti no tā paša avota, ja klasteru skaitlis ir izvēlēts pareizi. Paraugu tuvums klasteros iekšā var tikt

mērīts ar *p*-vērtību palīdzību, aprēķinātu atbilstošam Birojviesošanas *T*-kvadrāta statistiķim. Datu nepiederošie un klasteringa algoritma nepilnīgums var dot mazas *p*-vērtības. Tomēr autori pieņem, ka to empīriskais sadalījums ir vismazāk koncentrēts oriģinālā, ka klasteru skaits ir izvēlēts pareizi. Minētā procedūra var būt aptuveni aprakstīta kā šādas sadales ar koncentrācijas testa secīgu lietošanu paaudze.

Atslēgvārdi: klasterings, Birojviesošanas kritēriji, klastera validācija

**I. Noviks.** Nervu-muskuļu savienojuma veidošanās paredzēšana, *Modelling and New Technologies*, 14.sēj., Nr.4, 2010, 73.–76. lpp.

Pirmo pāris nedēļu laikā pēc piedzimšanas neironi, kas inervē muskuli 'sacenšas' cits ar citu, lai inervētu maksimālu skaitu muskuļu šķiedru. Šīs sacensības rezultāts ir, ka mazāk aktīvie neironi inervē vairāk muskuļu šķiedru nekā aktīvākie neironi. Tas tiek saukts par 'lieluma principu' un tas nozīmē, ka mazāk aktīvie neironi vinnē vairākās sacensībās. Tas ir pārsteidzoši, skatoties uz *izolētas* muskuļu šķiedras (i.e. viena kauja) salīdzinoša priekšrocību pret aktīvākiem neironiem, kas tiek arī pētīts konkrētajā rakstā. Pielietojot spēles teorijas pieeju, mēs skaidrojam, iepriekšējā rakstā [1], kā tas var notikt. Šajā rakstā autori parāda dažus paredzējumus, kas izriet no viena modeļa.

Atslēgvārdi: spēles teorija, lieluma princips, skaitļošanas neirozinātne

## COMPUTER MODELLING & NEW TECHNOLOGIES

## ISSN 1407-5806 & ISSN 1407-5814(on-line)

### **EDITORIAL BOARD:**

Prof. Igor Kabashkin (Chairman of the Board), *Transport & Telecommunication Institute, Latvia;*Prof. Yuri Shunin (Editor-in-Chief), *Information Systems Management Institute, Latvia;*Prof. Adolfas Baublys, *Vilnius Gediminas Technical University, Lithuania;*Dr. Brent D. Bowen, *Purdue University, USA;*Prof. Olgierd Dumbrajs, *Helsinki University of Technology, Finland;*Prof. Eugene Kopytov, *Transport & Telecommunication Institute, Latvia;*Prof. Arnold Kiv, *Ben-Gurion University of the Negev, Israel;*Prof. Juris Zakis, *University of Latvia;*Prof. Edmundas Zavadskas, *Vilnius Gediminas Technical University, Lithuania.*

Host Organization: Transport and Telecommunication Institute

### **Supporting Organizations:**

Latvian Transport Development and Education Association Latvian Academy of Sciences Latvian Operations Research Society

# THE JOURNAL IS DESIGNED FOR PUBLISHING PAPERS CONCERNING THE FOLLOWING FIELDS OF RESEARCH:

- mathematical and computer modelling
- mathematical methods in natural and engineering sciences
- physical and technical sciences
- computer sciences and technologies
- semiconductor electronics and semiconductor technologies
- aviation and aerospace technologies
- electronics and telecommunication
- navigation and radar systems
- telematics and information technologies
- transport and logistics
- economics and management
- social sciences

In journal articles can be presented in English. All articles are reviewed.

### EDITORIAL CORRESPONDENCE

Transporta un sakaru institūts (Transport and Telecommunication Institute) Lomonosova iela 1, LV-1019, Riga, Latvia. Phone: (+371) 67100593. Fax: (+371) 67100535. E-mail: journal@tsi.lv, www.tsi.lv

### COMPUTER MODELLING AND NEW TECHNOLOGIES, 2010, Vol. 14, No.4

Scientific and research journal of Transport and Telecommunication Institute (Riga, Latvia) The journal is being published since 1996.

### Computer Modelling & New Technologies \* Preparation of publication

## **PREPARATION OF CAMERA-READY TYPESCRIPT: COMPUTER MODELLING AND NEW TECHNOLOGIES**

- 1. In order to format your manuscript correctly, see the Page Layout Guideline for A4 (21 cm x 29,7 cm) paper size. Page Layout should be as follows: Top 3 cm, Bottom 3 cm, Left 3 cm, Right 3 cm.
- 2. Maximum length for the article is 10 pages.
- 3. Application of other Styles with the exception of Normal is impossible!
- 4. Articles should be Times New Roman typeface, single-spaced.
- 5. The article should include:
  - title;
  - author's name(s) and information (organisation, city, country, present address, phones, and e-mail addresses);
  - abstract (100–150 words);
  - keywords (max. about six);
  - introduction clearly explaining the nature of the problem, previous work, purpose and contribution of the research;
  - description of proper research;
  - conclusion section (this is mandatory) which should clearly indicate advantages, limitations and possible applications;
  - references.
  - **Attention!** First name, last name, the title of the article, abstract and keywords must be submitted in the English and Latvian languages (in Latvian it is only for Latvian authors) as well as in the language of the original (when an article is written in different language).
- 6. The text should be in clear, concise English (or other declared language). Please be consistent in punctuation, abbreviations, spelling (*British English*), headings and the style of referencing.
- 7. *The title of the article* 14 point, UPPERCASE, style Bold and centred.
- 8. Author's names centred, type size 12 point, Upper and lower case, style Bold Italic.
- 9. Author's information 10 point, Upper and lower case, style Italic, centred.
- 10. *Abstract and keywords* 8 point size, style Normal, alignment Justify.
- 11. *The first level Headings* 11 point, Upper and lower case, style Bold, alignment Left. Use one line space before the first level Heading and one line space after the first level Heading.
- 12. *The second level Headings* 10 point, Upper and lower case, style Bold, alignment Left. One line space should be used before the second level Heading and 1/2 line space after the second level Heading.
- 13. *The third level Headings* 10 point, Upper and lower case, style Italic, alignment Left. One line space should be used before the second level Heading and 1/2 line space after the third level Heading.
- 14. *Text* of the article 10 point, single-spaced, alignment Justify.
- 15. The set of *formulas* on application of fonts, signs and a way of design should be uniform throughout the text. The set of formulas is carried out with use of editors of formulas MS Equation 3.0 or MathType. The formula with a number the formula itself should be located on the left edge of the text, but a number on the right. Font sizes for equations are: 11pt full, 7pt subscripts/superscripts, 5pt sub-subscripts/superscripts, 16pt symbols, 11pt subsymbols.
- 16. All *Figures* must be centred. Figure number and caption always appear below the Figure, type size 8 point.

### Figure 1. This is figure caption

*Diagrams, Figures and Photographs* – must be of high quality, B in format \*.TIFF, \*.JPG, \*.BMP with resolution not less than 300 dpi. Also formats \*.CDR, \*.PSD are possible. Combination of Figures in format, for instance, \*.TIFF with elements of the in-built Figure Editor in MS Word is prohibited.

### Computer Modelling & New Technologies \* Preparation of publication

17. *Table Number and Title* – always appear above the Table. Alignment Left. Type size 8 point. Use one line space before the Table Title, one line space after the Table Title and 1/2 line space after the Table.

**Table 1.** This is an example of a Table

Heading	Heading	Heading
Text	Text	Text
Text	Text	Text

 References in the text should be indicated by a number in square brackets, e.g. [1]. References should be numbered in the order cited in the manuscript. The correct format for references is the following:

Article: author, title, journal (in italics), volume and issue number, year, inclusive pages
Example: 1. Amrahamsson, M., Wandel, S. A Model of Tearing in Third – Party Logistics with a Service Parts Distribution Case Study, *Transport Logistics*, Vol. 1, No 3, 1998, pp. 181-194.

*Book*: author, title (in Italics), location of publishers, publishers, year

Example: 2. Kayston, M. and W. R. Fried. *Avionic Navigation Systems*. New York: John Wiley and Sons Inc., 1969.

*Conference Proceedings:* author; title of an article; proceedings (in italics); title of a conference, date and place of a conference; publishing house, year, pages.

Example: 3. Canales Romero, J. A First Step to Consolidate the European Association of Aerospace Students in Latvia (Presented by the Munich Local Group). In: *Research and Technology – Step into the Future: Program and Abstracts. Research and Academic Conference, Riga, Latvia, April 7-11, 2003, Transport and Telecommunication Institute.* Riga: TTI, 2003, p. 20.

### 19. Authors Index

Editors form the author's index of a whole Volume. Thus, all contributors are expected to present personal colour photos with the short information on the education, scientific titles and activities.

### 20. Acknowledgements

Acknowledgements (if present) mention some specialists, grants and foundations connected with the presented paper. The first page of the contribution should start on page 1 (right-hand, upper, without computer page numbering). Please paginate the contributions, in the order in which they are to be published. Use simple pencil only.

21. Articles poorly produced or incorrectly formatted may not be included in the proceedings.