

РОССИЙСКАЯ АКАДЕМИЯ НАУК
ИНСТИТУТ ПРОБЛЕМ ИНФОРМАТИКИ

СИСТЕМЫ И СРЕДСТВА ИНФОРМАТИКИ

Специальный выпуск

МАТЕМАТИЧЕСКИЕ МОДЕЛИ
В ИНФОРМАЦИОННЫХ ТЕХНОЛОГИЯХ

Ежегодник основан в 1989 году

Главный редактор
член-корреспондент РАН И. А. СОКОЛОВ



Москва
ИПИ РАН
2008

РЕДАКЦИОННАЯ КОЛЛЕГИЯ:

Член-корреспондент РАН И. А. СОКОЛОВ (главный редактор),
доктор физико-математических наук С. Я. ШОРГИН
(заместитель главного редактора),
кандидат технических наук В. Н. ЗАХАРОВ,
доктор технических наук В. Д. ИЛЬИН,
доктор физико-математических наук Л. А. КАЛИНИЧЕНКО,
доктор технических наук В. А. КОЗМИДАДИ,
доктор технических наук К. К. КОЛИН,
доктор физико-математических наук А. В. ПЕЧИНКИН,
доктор геолого-минералогических наук Р. Б. СЕЙФУЛЬ-МУЛЮКОВ,
доктор технических наук И. Н. СЕНИЦЫН,
кандидат технических наук А. В. ФИЛИН,
кандидат физико-математических наук С. А. ХРИСТОЧЕВСКИЙ,
О. В. ЛОМАКИНА (ответственный секретарь)

Рецензент:

доктор физико-математических наук В. И. СЕНИЦИН

Системы и средства информатики: Спец. вып. Математические модели в информационных технологиях / Отв. ред. И. А. Соколов. — М.: ИПИ РАН, 2008. — 144 с. — ISBN 5-902030-36-6 (978-5-902030-36-2).

Специальный выпуск “Математические модели в информационных технологиях” ежегодника трудов Института проблем информатики РАН “Системы и средства информатики” посвящён одному из важнейших направлений фундаментальных исследований, проводимых в ИПИ РАН, — развитию математических методов, используемых в информатике и информационных технологиях.

В сборнике представлены работы сотрудников Института, относящиеся к следующим областям исследований: стохастические модели и основанные на них информационные технологии, модели информационных процессов и структур в информационно-вычислительных и телекоммуникационных системах, разработка новых математических методов, ориентированных на задачи информатики.

Для научных работников, инженеров, преподавателей, аспирантов, студентов вузов, интересующихся современным состоянием исследований в области математических методов информатики.

PREFACE

The present volume includes the articles devoted to the development and the study of mathematical and computer models intended to applied problems arising in various fields. Most of the presented researches are performed by the scientists from ORT Braude College (Karmiel, Israel) and the Institute of Informatics Problems of the Russian Academy of Sciences (Moscow, Russian Federation). The Joint Israeli-Russian Symposium on Stochastic models, which took place in Karmiel (2006), has initialized cooperation between these scientific centers. The XXVI International Seminar on Stability of the Stochastic Models (Naharya, 2007), organized under the aegis of the mentioned institutes, became an important stage in the development of this partnership.

The mutual agreement between two organizations, signed in 2006, suggests performing of joint projects and the specialists exchange in the research and the educational areas. The present volume is one of results of this fruitful cooperation.

Persistent contacts between two centers significantly influence the subjects and encourage the development of new fundamental and applied outcomes.

The presented volume also contains results obtained in cooperation of experts from two institutes with numerous researches from other leading scientific centers such as the Lomonosov Moscow State University, the Russian State Humanitarian University (Russia), the Technion-the Israeli Technological University, University of Haifa (Israel), the Norwegian University of Science and Technology, University of Salerno (Italy) etc.

The joint paper of the authors from the Institute of Informatics Problems of the Russian Academy of Sciences and the Holon Institute of Technology is written in the framework of mutual agreement signed in 2005.

Many results presented in the volume can be consider as an extension of the results presented on the Joint Israeli-Russian Symposium on Stochastic models (Karmiel, 2006) and the XXVI

International Seminar on Stability of the Stochastic Models (Naharya, 2007).

All the manuscripts have undergone reviewing process and have been selected by leading scientists of ORT Braude College and the Institute of Informatics Problems.

The Editors hope that the volume publication will encourage the cooperation between two institutes and will expand the fields of joint investigations performed by the scientists of Israel and Russia.

Editors:

Sergey Shorgin — Deputy Director of the Institute of Informatics Problems of the Russian Academy of Sciences

Zeev (Vladimir) Volkovich — Head of the Software Engineering Department, ORT Braude College, Israel

Предисловие

Настоящий специальный выпуск ежегодника трудов Института проблем информатики РАН «Системы и средства информатики» включает статьи, посвященные разработке и исследованию математических и компьютерных моделей, предназначенных для решения прикладных проблем в различных отраслях науки и практики. Представленные исследования выполнены в основном специалистами из Института проблем информатики Российской академии наук (Москва, Российская Федерация) и Колледжа ОРТ Брауде (Кармиель, Израиль). Начало сотрудничеству между этими двумя научными центрами было положено в ходе Израильско-Российского симпозиума по стохастическим моделям (Кармиель, 2006), а важным этапом в развитии совместных работ стал XXVI Международный семинар по устойчивости стохастических моделей, прошедший по эгидой названных организаций в Нахрии в 2007 году.

Соглашение о сотрудничестве между нашими организациями, подписанное в 2006 году, предусматривает осуществление совместных проектов и обмен специалистами в научно-исследовательской и образовательной областях. Настоящий сборник статей является одним из результатов этого плодотворного сотрудничества.

Постоянные научные контакты между специалистами обеих организаций оказывают большое положительное влияние на рас-

ширение тематики исследований и достижение принципиально новых фундаментальных и прикладных результатов.

В сборнике представлены также работы, содержащие результаты исследований ученых из других ведущих научных центров, таких как Московский государственный университет имени М.В.Ломоносова, Российский государственный гуманитарный университет (Россия), Технион-Израильский технологический университет, Хайфский университет (Израиль), а также Норвежский университет науки и технологии, Университет Салерно (Италия) и др., выполненных в рамках совместных проектов с учеными из Института проблем информатики РАН и Колледжа ОРТ Брауде.

Совместная статья авторов из Института проблем информатики РАН и Холонского Технологического Института выполнена в соответствии с Соглашением о сотрудничестве, заключенным между этими организациями в 2005 г.

Многие результаты, представленные в сборнике, явились дальнейшим развитием результатов, представленных в 2006 на Израильско-Российский симпозиуме по стохастическим моделям и в 2007 году на XXVI Международном семинаре по устойчивости стохастических моделей.

Статьи сборника прошли рецензирование и отбор ведущими учеными Института проблем информатики РАН и Колледжа ОРТ Брауде.

Составители сборника выражают надежду, что это издание станет придаст новый импульс сотрудничеству между нашими двумя организациями и расширит сферу совместных исследований и разработок ученых и специалистов Израиля и России.

Редакторы-составители:

Зеев (Владимир) Волькович — декан факультета инженерного программирования, Колледж ОРТ Брауде, Израиль

С. Я. Шоргин — заместитель директора Института проблем информатики Российской академии наук

ON A SIMULATION APPROACH TO CLUSTER STABILITY VALIDATION

Zeev Barzily

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

Mati Golani

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

Zeev Volkovich

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

In the current paper we outline a new approach to the “true number of clusters” determination problem. Our method combines both the stability and density concentration approaches. In the spirit of the density estimation methodology, we consider each cluster as an island of “high” density of items in a sea of “low” density. In addition, following the cluster steadiness concept, we suggest that these islands are “resistant” to a random noise. In other words, we believe that adding noise to the attributes of the data elements does not change the clusters structure. A second novelty of our approach is the proposition to measure the similarity between source-data clusters and noisy-data clusters by means of two sample test statistics, represented by probability metrics-distances. Such a pair seems as an appropriate database for the true number of clusters determination. As a consequence of the high resemblance between these samples, within the partitions, the similarity is expected to be amplified under the true number of clusters. According to our model, the true number of clusters corresponds to the empirical distance distribution which is most concentrated at zero. Thus, our procedure can be considered as the creation of an empirical normalized distance distribution, followed by testing its concentration at zero. This test is carried out by means of the sample mean and the size of the sample first quartile.

1. Introduction

Clustering is a technique for intellectual data analysis. It is applied in disciplines like social sciences, biology and computer science, in attempts to acquire intuitive understanding of the data meaning. Clustering of data items is a basic tool for achieving

this goal. Group membership is typically discovered via an iterative clustering procedure which employs a distance-like function that measures the resemblance between data points. The result of the clustering process is a partition of the source data which is characterized by the highest quality score. Excluding the data itself, two crucial input parameters are required for the application of iterative clustering procedures: an initial partition and a suggested number of clusters. In many practical applications, the problem of choosing the right number of clusters is still, more or less, unsolved. It is well known that this important task is “ill posed” [Jain and Dubes, 1988], [Gordon, 1999]. For example, the ‘correct’ number of clusters can depend on the scale in which the data is presented (see, for example [Chakravarthy and Ghosh, 1996.]). Many approaches have been offered to solve this problem. Up to now, none of them has been recognized as superior to the others. Most methods can be classified into the following groups:

- Multivariate statistical indexes which compare dispersions within and between the clusters;
- Stability (similarity, merit) functions that evaluate the consistency of labels assignments to sample elements;
- Density estimation approaches.

The following papers are members of the first group: [Dunn74], [Hubert74], [Calinski74], [Hartigan85], [Krzanowski85], [Sugar03] and [Tibshirani01]. In papers belonging to the second group, stability is understood as the fraction of times items, or pairs of items, are assigned to the same cluster. These methods follow the basic idea that, a clustering algorithm, repeatedly applied to random samples from a population, has to construct similar clusterings (see, for example, [Roth02], [Levine01] and [Ben-Hur02]). Clustering methods, based on estimation of the underlying data density, assume that the clusters are related to modes of the probability density function. These clustering procedures allocate an item to one “domain of attraction” of a density peak [Wishart69], ([Hartigan75] Section 11) and [Hartigan81]. The number of clusters is identified here, as the quantity of disjoint intervals having densities exceeding a predefined value (see, for example [Cuevas00], [Cuevas01] and [Stuetzle03]).

In the current paper we address to the “true number of clusters” determination problem. Our method combines the stability and the density concentration approaches. Specifically, in the spirit of the density estimation methodology, we consider each cluster as an

island of “high” density of items in a sea of “low” density. In addition, following the cluster steadiness concept, we suggest that these islands are “resistant” to a random noise. In other words, we presume that adding noise to the attributes of the data elements does not affect much the clusters structure if the number of clusters is chosen correctly. Another novelty of our approach is that, we propose to measure the similarity between source-data clusters and noisy-data clusters by means of two sample test statistics, represented by probability metrics-distances. Such a pair seems appropriate for the true number of clusters determination because high resemblance between these samples, within the partitions, is expected to be present only under the true number of clusters.

The method proposed here is outlined as follows: Pairs of samples are considered for each tested number of clusters such that the first one is drawn without replacement from the data. The second one is obtained by adding a random noise to the elements of the first sample. The noise is simulated as a sequence of independent identical distributed random variables. The distances between the pairs are measured via a probability metric between the samples within the clusters. This kind of metric appears in two sample tests where it is required to settle on whether two specified samples are derived from the same population. In this paper we use the probabilistic metrics which compare the mean kernel distance of the pooled clustered sample to the mean kernel distances obtained in both drawn clustered samples. In order to employ these distances each of the samples is clustered twice: alone and together with the other. Actually, distances compare the partitions found. Once the obtained metrics values are sufficiently small, the partitions are considered close. Samples outliers and drawbacks of clustering algorithms contribute to the instability of results obtained. This obstacle is overcome by repeating the described routine many times with appropriate outcomes normalization. According to our model, the true number of clusters corresponds to the empirical distance distribution which is most concentrated at zero. Hence, our procedure can be considered as the creation of an empirical normalized distance distribution, and afterwards the testing of its concentration at zero. This test can be carried out by means of numerous simple statistics such as the sample mean and the size of the sample first quartile. We choose the “true” number of clusters as the one, which minimizes the normalized average probability distance or the normalized sample first quartile.

2. The approach

We concern a finite subset $X = \{x_1, x_2, \dots, x_n\}$ of the d -dimensional Euclidean space \mathbb{R}^d . For a specified subset $S \subset X$ a partition $\Pi_k(S)$ of S is a family of subsets of S

$$\Pi_k(S) = \{\pi_1(S), \pi_2(S), \dots, \pi_k(S)\},$$

for which the following holds:

$$\bigcup_{j=1}^k \pi_j(S) = S$$

and

$$\pi_i(S) \cap \pi_j(S) = \emptyset, \quad i \neq j.$$

The elements of $\Pi_k(S)$ are called *clusters*. We assume that a clustering algorithm Cl is available. The algorithm has the sample and the given number of clusters k as input parameters. A data partition is the algorithm output.

In the framework of our approach we consider for each sample S , drawn from X , its noisy version \tilde{S} composed from the elements $\tilde{x}_i = x_i + \varepsilon_i$ where $x_i \in S$ and ε_i are independent identically distributed random variables, $i = 1, \dots, n$. We further assume, a very common point of view, that the random variables ε_i are normally distributed having the mean zero, $\varepsilon_i \sim N(0, \sigma)$, where the standard deviation σ presents the noise altitude. Certainly, σ can be considered a control parameter. Very small values of σ can lead to the procedure instability because, in this case, both samples are actually drawn from the same population. Big values of σ , actually yield noisy samples which are uniformly distributed. This situation is, also, not desirable. Effects of the σ size for three cases corresponding to $\sigma = 0.1, 0.3, 0.5$, for a sample drawn from a two-cluster population, are presented on Fig. 1 (see supplementary sheet 1).

In order to characterize a resemblance between two given samples, say S_1 and S_2 , having a size n , we use the kernel R -distances presented in [Klebanov03] and [Klebanov05]. This distance depends on a real symmetric negative definite kernel $K(x, y)$. Such a distance can be defined as:

$$d^2(S_1, S_2) = 2\Delta(S_1, S_2) - \Delta(S_1, S_1) - \Delta(S_2, S_2),$$

where

$$\Delta(S_1, S_2) = \frac{1}{n^2} \sum_{x_i \in S_1} \sum_{y_j \in S_2} K(x_i, y_j).$$

Similar distances have been suggested in [Baringhaus04] and [Zech05]. In our methodology we use these distances as follows: For given number of clusters k and given noise level σ , we construct M pairs of samples $S_1^{(m)}$ and $S_2^{(m)}$ ($m = 1, \dots, M$) such that $S_1^{(m)}$ is selected randomly, without replacement, from X . The sample $S_2^{(m)} = \widetilde{S_1^{(m)}}$ is its noisy version $S_1^{(m)}$ simulated each time independently according to the described above procedure. We introduce

$$S^{(m)} = S_1^{(m)} \cup S_2^{(m)}$$

and its partition

$$\Pi_k^{(m)} = Cl(S^{(m)}, k).$$

Now we consider

$$S_{1l}^{(m)} = S_1^{(m)} \cap \pi_l^{(m)}(S^{(m)}), \quad S_{2l}^{(m)} = S_2^{(m)} \cap \pi_l^{(m)}(S^{(m)})$$

as the subsets of $S_1^{(m)}$ and $S_2^{(m)}$, respectively, the elements of which are members of the cluster $\pi_l^{(m)}(S^{(m)})$. We quantify dissimilarities $D_l^{(m)}$, $l = 1, \dots, k$ amid these sets inside $\pi_l^{(m)}(S^{(m)})$ with the help of the R -distances: $D_l^{(m)} = d^2(S_{1l}^{(m)}, S_{2l}^{(m)})$. Thus, the dissimilarity between $S_1^{(m)}$ and $S_2^{(m)}$ is defined as

$$D^{(m)}(S_1^{(m)}, S_2^{(m)}) = \frac{1}{k} \sum_{l=1}^k \text{abs} \left(D_l^{(m)} \right).$$

The outline of the suggested algorithm is as follows:

- k^* maximal number of clusters to be tested;
 - n size of the samples;
 - M number of samples pairs;
 - T number of the averaged distance values;
 - K a negative definite kernel;
 - σ a noise level.
1. For $k = 2 : k^*$ do
 2. For $t = 1 : T$ do

-
3. For $m = 1 : M$ do
 4. $S_1^{(m)} = \text{sample}(X, n)$;
 5. Simulate a noisy sample $S_2^{(m)}$;
 6. $dis^{(m)} = D^{(m)}(S_1^{(m)}, S_2^{(m)})$;
 7. end for 3
 8. $Dis_t = \frac{1}{M} \sum_{m=1}^M dis^{(m)}$;
 9. end for 2
 10. Normalize the array Dis ;
 11. end for 1
 12. The selected number of clusters is the one that corresponds to the most concentrated distribution of Dis at zero.

3. Experiments

In order to evaluate the functioning of our method, we present several numerical experiments on synthetic and real datasets. The drawn samples are clustered by the standard k -means algorithm. The concentration of the empirical distances distributions is characterized by means of the average and the C_{25} (the 25th percentile). The normalization is provided by the sample 95th percentile. The Gaussian kernel

$$K(x, y) = \|x - y\|^2$$

is used. We perform 10 trials for each experiment with $k^* = 7$, $n = 200$, $M = T = 30$ and $\sigma = 0.1$. The results are presented via the error-bar plots of the average and the C_{25} of the distances in the trials. The sizes of the error bars are two standard deviations.

3.1. Synthetic datasets. The considered synthetic datasets have been simulated as mixtures of NC , two-dimensional, Gaussian distributions having the same standard deviation denoted by SD .

The components means are situated on the unit circle with an equal distance between each neighboring pair. Each component contains 3000 points. We denote such datasets as G_NC_SD . A scatter-plot of such a dataset is presented in the following figure.

As we can see, in the presented dataset, the components are not disjoint. However the proposed succeeds to point out the true number of clusters.

Another dataset considered is a G_4_04 dataset. Obviously, this collection is less separable. Consequently, the outcomes here are vaguer.

3.2 Real datasets. The first dataset is a text collection data which has been chosen from: <ftp://ftp.cs.cornell.edu/pub/smart/>. It includes

- DC0–Medlars Collection (1033 medical abstracts);
- DC1–CISI Collection (1460 information science abstracts);
- DC2–Cranfield Collection (1400 aerodynamics abstracts).

This dataset was considered in several papers (see, for example [Kogan03]. Applying the well-known “bag of words” approach, 600 “best” terms were selected (see, for example, [Dhillon03] for term selection details). The dataset is embedded into Euclidian spaces having dimensions of 600. A dimension reduction is provided by the Principal Component Analysis. The data is recognized to be well separated by means of the two leading principal components. Thus, we use this data representation in our experiments. Accordingly, we operate with this data representation in our experiments. As demonstrated in Fig. 5 (see supplementary sheet 1), our method precisely determines the true number of clusters in this case.

The second considered dataset is the, well known, Iris Flower Dataset available at <http://fmwww.bc.edu/ec-p/data/micro/iris.dta>. This four dimensional dataset includes information on three flowers types: 0 — Iris Setosa, 1 — Iris Versicolour and 2 — Iris Virginica. 50 elements are available for each flower type. It is well known that one of the sets is linearly separable from the others while the other two are not. Analysis of this dataset encountered difficulties in detecting all three clusters, (see for example [Roth02]) where two clusters were detected. These difficulties are due to the nonlinear separation of the third cluster. The fact has been explained such that the third cluster had not been detected due to the nonlinearity of its separation from others. In the paper [Roth02] a variant of the k -means algorithm has been used. Such an algorithm produces

linear boards between the clusters. We also use a variant of the k -means algorithm; however our methodology detects a three cluster structure.

4. Summary

We propose a new methodology for cluster stability. It combines the stability and density concentration approaches. In the spirit of the density estimation point of view, we consider each cluster as an island of “high” density in a sea of “low” density. We compare pairs of samples, where the second sample is obtained from the first by adding noise to each of its elements. We find that adding an appropriate noise to the attributes of the data elements does not have adverse effect on the detected clusters structure. We offer to measure dissimilarities between a clustered sample and its clustered noisy version by means of kernel two sample test statistics. Several provided experiments demonstrate that the cluster islands are detected.

According to our model, the true number of clusters corresponds to the empirical distance distribution which is most concentrated at zero. Our procedure can be considered as the creation of an empirical normalized distance distribution, followed by testing its concentration at zero. This test is carried out by means of the sample mean and the size of the sample first quartile

5. Future work

Several options of extended work are considered:

- Noise stability limit — applying the algorithm on varying levels of noisy data sets in order to determine the noise level that causes the algorithm to collapse.
- Creation of a theoretical model which allows estimating the optimal noise altitude.
- A precision/recall comparison with known algorithms in order to test performance of the algorithm.

References

- [Baringhaus04] *Baringhaus L. and Franz C.* On a new multivariate two-sample test // *Journal of Multivariate Analysis.* 2004. V. 88, No. 1. P. 190–206.

-
- [Ben-Hur02] *Ben-Hur A. et al.* A stability based method for discovering structure in clustered data // In: Pacific Symposium on Biocomputing. 2002. P. 6–17.
- [Calinski74] *Calinski R. and Harabasz J.* A dendrite method for cluster analysis // Commun Statistics. 1974. V. 3, No. 1. P. 27.
- [Chakravarthy96] *Chakravarthy S. V. and Ghosh J.* Scale-based clustering using the radial basis function network // In: IEEE Transactions on Neural Networks. 1996. V. 7, No. 5. P. 1250–1261.
- [Cuevas00] *Cuevas A. et al.* Estimating the number of clusters // The Canadian Journal of Statistics. 2000. V. 28, No. 2. P. 367–382.
- [Cuevas01] *Cuevas A. et al.* Cluster analysis: A further approach based on density estimation // Computational Statistics and Data Analysis. 2001. V. 28. P. 441–459.
- [Dhillon03] *Dhillon I. et al.* Feature selection and document clustering // In: M. Berry, editor, A Comprehensive Survey of Text Mining. 2003. P. 73–100. Springer, Berlin, Heidelberg, N. Y.
- [Dunn74] *Dunn J. C. et al.* Well Separated Clusters and Optimal Fuzzy Partitions // J. Cybern. 1974. V. 4. P. 95–104.
- [Gordon99] *Gordon A. D.* Classification. — Chapman and Hall, CRC, Boca Raton, FL, 1999.
- [Hartigan75] *Hartigan J.* Clustering Algorithms. — John Wiley, N. Y., USA, 1975.
- [Hartigan81] *Hartigan J.* Consistency of single linkage for high-density clusters // Journal of the American Statistical Association. 1981. V. 76 P. 388–394.
- [Hartigan85] *Hartigan J.* Statistical theory in clustering // J. Classification. 1985. V. 2. P. 63–76.
- [Hubert74] *Hubert L. and Schultz J.* Quadratic assignment as a general data-analysis strategy // Br. J. Math. Statist. Psychol. 1974. V. 76. P. 190–241.
- [Jain88] *Jain A. and Dubes R.* Algorithms for Clustering Data. — Englewood Cliffs, Prentice-Hall, New Jersey, USA, 1988.
- [Klebanov03] *Klebanov L.* One class of distribution free multivariate tests // SPb. Math. Society, Preprint. 2003. V. 03, St-Petersburg, Russia (in Russian).
- [Klebanov05] *Klebanov L.* *N*-distances and their Applications. — Charsel University in Prague, The Karolinum Press, Prague, Czech, 2005.
- [Krzanowski85] *Krzanowski W. and Lai Y.* A criterion for determining the number of groups in a dataset using sum of squares clustering // Biometrics. 1985. V. 44. P. 23–34.

-
- [Kogan03] *Kogan J. et al.* Text mining with information-theoretical clustering // Computing in Science & Engineering. 2003. P. 52–59, November/December.
- [Levine01] *Levine E. and Domany E.* Resampling method for unsupervised estimation of cluster validity // Neural Computation. 2001. V. 13. P. 2573–2593.
- [Roth02] *Roth V. et al.* A resampling approach to cluster validation // In: COMPSTAT, available at <http://www.cs.uni-bonn.de/~braunm>, 2002.
- [Sugar03] *Sugar C. and James G.* Finding the number of clusters in a data set: An information theoretic approach // J. of the American Statistical Association. 2003. V. 98. P. 750–763.
- [Stuetzle03] *Stuetzle W.* Estimating the cluster tree of a density by analyzing the minimal spanning tree of a sample // J. Classification. 2003. V. 20, No. 5. P. 25–47.
- [Tibshirani01] *Tibshirani R. et al.* Estimating the number of clusters via the Gap statistic // J. Royal Statist. Soc. B. 2001. V. 63, No. 2. P. 411–423.
- [Tibshirani04] *Roth V. et al.* Stability-based validation of clustering solutions // Neural Computation. 2004. V. 16, No. 6. P. 1299–1323.
- [Wishart69] *Wishart D.* Mode analysis: A generalization of nearest neighbor which reduces chaining effects // In: Numerical Taxonomy. 1969. P. 282–311.
- [Zech05] *Zech G. and Aslan B.* New test for the multivariate two-sample problem based on the concept of minimum energy // The Journal of Statistical Computation and Simulation. 2005. V. 75, No. 2. P. 109–119.

AN APPROACH TO CLASSIFICATION OF COMPUTER SYSTEMS FAULTS LOCALIZATION MODELS

Sergey Frenkel

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

Eugene Levner

Holon Institute of Technology, Holon, Israel

Victor Zakharov

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

This paper suggests a characterization of fault localization strategies and a classification of probability-based search algorithms, which are widely used in the current practice of complex computer systems and networks systems maintenance and service. The characterization is performed in terms of some tuples including: characteristics of fault models, localization procedures, cost functions, and other factors. Such characterization (and an induced classification) can be used for a rational choice of search algorithms at early system design stages, for development of strategies of possible fault localization during the target systems maintenance and service. Search algorithms classification is based on the special notation like those used in the queuing theory.

1. Introduction

In computer systems development and service area, the important part of the system service cost falls into faults search ("fault management" [1, 2]), which is aimed to detect, diagnose and correct the possible faults during the system operations.

In order to provide a qualified fault diagnostic and localization for a computer system, it is necessary to involve numerous knowledge about all aspects of the faults search activity, from the search algorithms complexity to the hardware testability design issues.

The fault diagnostic and localization process in computer systems is a probabilistic in nature, that is the number of steps until any fault location can be recognized is a priori uncertainly. Fault search and localization probabilistic models are just constructed to

model various types of the uncertainties [2,3]. These models are, in fact, various search strategies, which specify the order in which probes and measurements required are to be performed, the type measurements (e.g., active probe selection or passive symptoms analysis [4]), etc.), and numerical models (e.g., Markov Chains Solver, Bayesian networks) [5, 6, 8]) used in the fault diagnostic planning.

However, in current computer engineering activity many conceptual aspects of the faults search problem mentioned above is often not taken into account. For example, in [4] the algorithms of finding minimal probe set is co-ordinated rather weakly with the problem of fault localization, in spite of the semantic closeness of these tasks. In [7] the test scheduling problem was considered without checkpoints choice. In other words, we deal with the problem how to use some information, obtained from the search algorithm analysis at early stages of the systems design to plan the failure detection, localization and correction in large size computer systems during their service. On the other hand, the lack of enough integrated view on the faults search process may prevent from understanding these subjects for students studying computer engineering.

The aim of this paper is to establish and describe possible relationships between various aspects of diagnostic modeling, extending our research of [17]. For this aim we classify well-known search algorithms from the viewpoint of their helpfulness for the problem of fault diagnostic and localization in complex computer systems. We determine the algorithms features, which can be useful in predicting the search cost by analyzing the relationship between those features, and features of the search space of an instance. As a result, we describe and ordering some properties of some modifications of well-known search algorithms using a classification, which is based on a special notation like that used in the queuing theory.

2. An informal description of faults search process

Let us consider an informal definition of the fault localization and diagnostic problem as a tuple $\langle S, F, P, Q \rangle$, where S is a description of target system (which should be developed) on a given step of design process (the level of system representation hierarchy, e.g., architectural description, or component-based representation), F is a set of possible faults in the system to be manufactured ("an object

of searching”, in fact), P is a plan (“strategy”) of possible faults search in the design during its functioning (actual or potential, in a working or in a testing modes), Q is a quality criterion (cost, time, accuracy, etc). The aim of this problem solution is to find a strategy of the faults localization at a given level of design hierarchy, and to estimate (predict) the cost of the fault localization activity as a part of the system service.

Let us call a detailed description of all considerable aspects of fault diagnostic process as a “general fault diagnostic model”, \mathbf{M}_{FD} , which has to expose information about the system structure that is useful to fault localization. Majority of algorithms of fault search in a computer system are based on the following knowledge:

- information about structure of a system diagnosed,
- information about possible faults properties which can affect the decision making about the faults discovery in the testing process,
- permissible testing strategy, including all taking into account resources limitations,
- information about the faults testing conditions, including possible errors (“noise”, uncertainties) during fault testing/probes.

We represent \mathbf{M}_{FD} as a tuple, each field of which belongs to corresponding domain relating to fault diagnostic and localization activity. In its part, each of fields may depend on some attributes of the diagnostic activity. At the current stage of thought, we propose the following parameterization of the model:

$$\mathbf{M}_{FD} = \langle \mathbf{SDF}, \mathbf{SPC}, \mathbf{FM}, \mathbf{FOM}, \mathbf{NHL}, \mathbf{SST} \rangle.$$

Let us consider in detail all these attributes.

SDF

The attribute **SDF** (“System Diagnostic Features”) describes technical features of the system which is under consideration with respect to fault diagnostic. From this point of view, the systems can be either fitted to active probing (e.g., ping or trace route command, an e-mail message, or a Web-page access request) or to passive testing. The major difference [4] is that we use an active probing approach versus a “passive” analysis of symptoms. So, we can consider three types of diagnostic systems:

- PFA is the abbreviation of “passive free access” to all structural elements (modules, components) of the system to test

them. In general, we may consider as an “element” a module of software, a module of hardware, nodes of a telecommunication networks,

- PFAC is the abbreviation of “passive free access to some set of checkpoints”, when any decisions about state of any inner module should be made only on the base of the state of a set of the checkpoints.
- AP is an active probing.

Thus, **SDF** can take symbolic values **SDF** = PFA, or **SDF** = PFAC, or **SDF** = AP.

SPC

The attribute **SPC** (“search precedence constraints”) deals with a structure of a search plan, which is determined considerably by the precedence constraints. For example, if we diagnose a D-flip-flop based structure, all test data must follow the “Set” or “Reset” signals.

The precedence relation constraints can be represented [7] by a directed precedence graph $G = (S, R)$, where the nodes in S correspond to the tests/probes (which, in fact, are considered in [7] as some “jobs” of the schedule theory), and an arc $(i; j) \in R$ corresponds to the precedence constraint, denoted $i \leq j$, that job i must precede job j . Thus, **SPC** may take values “0” or “1”, thereby **SPC** = 0 if there are no any precedence constraints in a search plan considered. In fact, **SPC** = 0 means that modules are inspected sequentially, whereas in other cases, for example, a number of modules can be expected simultaneously.

Note, that the precedence constraints can be used for adaptive tests building that are automatically tailored to a given level of a target system design [9].

FM

The attribute **FM** (“Fault model”) corresponds to a formal description of causes of errors which is considered as a discrepancy between an observed (or measured) value in a given place of the system and a true, specified, theoretically-expected correct value. There exist different ways to classify possible faults in computer systems. For example, the faults may be classified according to their duration time as: (1) permanent, (2) intermittent, and (3) transient. We will consider only permanent faults, as the current theory of transient faults search has still not taken shape [10]. In this case, the only explicit characteristic of the fault model, having

an impact on the search process complexity is the supposed number of possible faults in a specimen of the target system, that is either single (S) or multiple (M). So, the attribute FM may take the symbolic values S or M.

FOM

The attribute **FOM** (“Fault observation Model”) describes different possible ways of a fault observation. Usually faults are not directly observable and faults diagnosis is the process of locating the faults based on current observations (symptoms) and possibly further observations. The following fault detection conditions are considerable: Perf, that is “noise-free”[4] (or perfect) observation, and Imp, that is “noisy” (“imperfect”) [15]. In Perf any module recognized as a faulty one has its fault actually, and the module is fault-free actually if it is recognized as a fault free. In general, the Imp observations error may belong to both types: GF-error that is “good \rightarrow fail”, and FG that is “fail \rightarrow good” for another case of erroneous inspection of the modules. They can be both because of erroneous of modules inspection and uncertainty concerning the faults localization through symptoms observed. For example, the fail \rightarrow good error is mostly a result of not 100%-fault coverage.

So, the **FOM** attribute may take the symbolic values **FOM** = Perf, or **FOM** = Imp.

NHL

The attribute **NHL** (“Number of hierarchy levels”) characterizes a hierarchy of the fault localization process by the number of hierarchy levels L, where L is an integer value. It is possible to define several levels of hierarchy of the faults search problem as “level 1”, ..., “level n” (e.g., from a subsystem to a small system’s component, each of them can be a faulty with a fault model depending on the level). For example, in [7] two levels of sequential search for fault detection and diagnosis are considered, thereby each of level corresponding to a phase of the search process. The first phase deals with its domain for evidence of anomalous behavior. If evidence of such behavior is detected in a sub-domain, the second phase starts in which the target system element(s) in the sub-domain are tested to isolate the specific fault.

SST

SST (“search stopping criterion”) may take the following values:

- FFD — up to the first fault detection,
- AFD — up to all fault detection,

Table 1

The attributes of the fault diagnostic model

Attr.	Range of values	Sense of values	Destination of attribute
SDF	PFAC, PFA, AP	PFAC means passive free access to some points, PFA is passive free access to all modules, AP is active one	System diagnostic features
SPC	0, 1	0, if the modules are inspected sequentially, precedence constraints in a search plan considered, 1 otherwise	Search precedence constraints
FM	S, M	S is a single fault model, M is a multiple one	Fault model
FOM	Perf, Imp	“Perf” means perfect (noise-free) observation, Imp is imperfect (noisy) observations	Fault observation model
NHL	Any integer	The number of defined hierarchy levels for a given faults search problem	Number of hierarchy levels
SST	FFD, AFD, RS	FFD means a search up to the first fault detection, AFD means a search up to all fault detection, RS means up to a threshold value of resources.	Search stopping criterion

— RS — up to a threshold value of resources (e.g., number of steps, time restriction).

All considered attributes are listed in the Table 1.

Note, that among the attributes of the general diagnostic model also some attributes of “Cost model” should be used, that is a cost value that determines an optimal order in which the components should be tested. It may be a deterministic or random. Certainly, the type of a cost function (or “loss function” [11, 12, 13]) also affects the optimal ordering, but this influence takes place rather on a quantitative characteristic of the problem solving than on the structure of search model (Section 4). However, because of the restricted place, we will not consider models with the random cost functions [7] in this paper.

3. Search algorithm classification

There are many ways to classify fault localization algorithms, but the main thing is that possible algorithms with some useful

features should be chosen from the classes. We can use Classification of Queuing models as a pattern for our classification.

3.1. Kendall Classification of Queuing Systems. *Queuing Theory (QT)* is a collection of mathematical models of various queuing systems that take as inputs parameters of the above elements and that provide quantitative parameters describing the system performance.

The Kendall classification of queuing systems [18] exists in several modifications. The most comprehensive classification uses 6 symbols:

$$\mathbf{A}/\mathbf{B}/s/q/c/p$$

where:

A is the arrival pattern (distribution of intervals between arrivals).

B is the service pattern (distribution of service duration).

s is the number of servers.

q is the queuing discipline (FIFO, LIFO, ...). Omitted for FIFO or if not specified.

c is the system capacity. Omitted for unlimited queues.

p is the population size (number of possible customers). Omitted for open systems.

These symbols are used for arrival and service patterns:

M is the Poisson (Markovian) process with exponential distribution of intervals or service duration respectively.

E_m is the Erlang distribution of intervals or service duration.

D is the symbol for deterministic (known) arrivals and constant service duration.

G is a general (any) distribution.

GI is a general (any) distribution with independent random values.

Examples:

D/M/1 = Deterministic (known) input, one exponential server, one unlimited FIFO or unspecified queue (test sequence), unlimited customer population.

M/G/3/20 = Poisson input, three servers with any distribution, maximum number of customers 20, unlimited customer population.

D/M/1/LIFO/10/50 = Deterministic arrivals, one exponential server, queue is a stack of the maximum size 9, total number of customers 50.

Each sequence of symbols mentioned above represents some applied features of a specific system classified. Let us consider a possibility to extend this approach to classification of search algorithms used in the fault management.

3.2. QT-like notation for search algorithms classification.

Following this notation, we may suggest to describe and present different fault search problems from the area of computer system diagnosis by using a 5-attribute notation $(\alpha|\beta|\gamma|\delta|\varepsilon)$ in more economical form.

Table 2

Search algorithms classification symbols

Attr.	Notation and values	Meaning	Purpose
α	Active/Passive probing	(see above)	System type of the probing
β	Fault model: s, m, b, p, i, t	(see above)	Type of the failure
γ	Character of input information	(see above)	Type of the observed data
δ	Type of inspections: perf/imperf; prec	perf — perfect, imperf — imperfect; prec — there are precedence relations between tests	Type of inspections
ε	Search scenario: Be, M, Ba, P, h; ...	Be-Bernoulli, M-Markov, Ba-Bayes, P-Poincarre, h-hierarchical, etc.; Min-loss, Min-cost, Max-reward, etc.	Type of the search scenario and criteria

In what follows, we will consider two models:

$$(\alpha|\beta|\gamma|\delta|\varepsilon) = (PP|s|pr|imperf\text{ Be, Min-cost}),$$

and

$$(PP|s|pr|imperf|, M, \text{Max-reward}).$$

Fast real-time algorithms for the optimal search for the hidden faults in computer systems are developed [12].

They are based on the following a priori information, obtained in advance using expert's estimations:

- The probability that the i th specific module is failed,
- The probability (risk) of an unsuccessful search (“over-looking”),
- The expected cost and time of search trials for each individual module.

The search procedure uses the concept of the *dynamic effectiveness* of each trial, which is strictly defined in this paper and depends on time, cost and risk characteristics of the module, as well as on the search stage.

The necessary and sufficient conditions for the search optimality are found which claim that the linear, exponential and logarithm utility functions, and only these functions, guarantee that the *local* real-time search procedure provides the *global optimum*.

Model 1. Min-Cost Search: Costs depending on searching time

Consider collection \mathbf{K} of $|\mathbf{K}| = N$ independent stochastically failing modules and assume that \mathbf{K} contains only one failed module which we are to find with a minimum cost.

Assume that the cost is accumulated from the moment the search starts till the failed module is discovered and the failure is eliminated. Assume that a series of *sequential tests* must be performed for different modules, and the goal of the search is to minimize the expected cost of the search before the failure is localized and eliminated. We assume that the search is *imperfect*, which means that there is a risk that the test does not extract the failed module when it is examined, that leads to the situation when the same module can be examined several times. A search sequence, generally speaking, may become infinitely long. We assume that each module i is characterized by the following parameters:

- Prior probability p_i that i is failed;
- Prior probability a_i that i will not be extracted, by mistake, when the search test examines it,
- Expected time t_i to inspect module i ,
- Cost coefficients:
 - “cost rate” c_i , per unit time (for linear costs),
 - “maximal acceptable” C_i (for exponentially growing costs).

The optimal search strategy depends on the cost functions and the search scenario.

In Model 1 we consider the following scenario:

- The modules are inspected sequentially;
- For any specified search sequence and any location of the failed element, the outcomes of the searches are independent.
- The search is *finished* when the failed module is *found and the failure is eliminated*.

Each sequential inspection strategy specifies an infinite sequence s of inspected modules:

$$s = (s[1], s[2], \dots, s[n], \dots)$$

which states that at the n th step of s the module $s[n]$ is inspected.

For any s we introduce the following notation:

- $P[n, s]$: probability that the failure is detected at the n th step of s ;
- $t[d, s]$: time elapsed till the moment when the module $s[d]$ is examined at step d of sequence s ;
- $c[d, s]$: the cost rate assigned to module $s[d]$ in s ;
- $C[d, s]$: the maximal possible cost assigned to module $s[d]$ in s .

Then the time $T[n, s]$ spent to detect and eliminate the failure at the n th step of s will be:

$$T[n, s] = \sum_{d=1, \dots, n} t[d, s].$$

- The linear cost for n th step of s will be: $c[n, s]T[n, s]$;
- The exponentially growing cost will be:

$$C[n, s](1 - \exp(-dT[n, s])).$$

Then the total expected linear cost of sequence s will be:

$$L(s) = \sum_{d=1, \dots, \inf} P[n, s]c[n, s]T[n, s].$$

- The expected exponential costs will be:

$$E(s) = \sum_{d=1, \dots, \inf} P[n, s]C[n, s][1 - \exp(-dT[n, s])].$$

Let us define parameters Q_{ij} and R_{ij} :

$$Q_{ij} = p_i(a_i)^{(j-1)}(1 - a_i);$$

$$R_{ij} = C_i Q_{ij} \exp(-dt_i) / [1 - \exp(-dt_i)].$$

Theorem 1. *Let the values of ratio $c_{ij}Q_{ij}/t_i$ for all i and j be arranged in non-increasing order, and let s^* be such that if the ratio $c_{ij}Q_{ij}/t_i$ is the n th largest one in the ordering then the n th step of s^* is the j th search of module i . Then s^* is optimal.*

Theorem 2. *Let the values of ratio R_{ij} for all i and j be arranged in non-increasing order, and let s^* be such that if the ratio R_{ij} is the n th largest one in the ordering then the n th step of s^* is the j th search of module i . Then s^* is optimal.*

Parameters $c_{ij}Q_{ij}/t_i$ and R_{ij} are called *the dynamic effectiveness*, or *attractiveness* of the module i at its j th search.

Example: $Q_{11} < Q_{21} < Q_{31} < Q_{12} < Q_{13} < Q_{14}$.

Model 2. Max-Reward Search with Learning:

The faster failure is detected the larger reward.

We assume that each module i is characterized by the following parameters:

- Prior probability p_i that this module i is failed;
- Prior probability a_i that i will not be extracted, by mistake, when the search test examines it,
- Expected time t_i to inspect i ,
- Reward parameters:
 - a reward W_i offered at time 0 for finding the needed info,
 - “a discount factor” d_i , showing the decrease of reward over time,
 - “a learning factor” b_i , $0 < b_i < 1$, (see Sweat 1979) describing a relative increase of the probability p_i of finding the needed info: the latter probability contains factor $\prod_{i=1, \dots, N} (b_i)^{u(i)}$, where $u(i)$ is the number of times module i has been inspected before.

The optimal search strategy will depend on the reward functions and the search scenario.

In this model we consider the following scenario:

- The modules are inspected sequentially;
- For any specified search sequence and any location of the failure, the outcomes of the searches are dependent on the previous searches;

- The search is *finished* when the failure is localized and eliminated;
- Each sequential inspection strategy specifies an infinite sequence s of modules:

$$s = (s[1], s[2], \dots, s[n], \dots)$$

which states at the n th step of s the module $s[n]$ is inspected.

For any s we introduce the following notation:

- $W[n, s]$: the reward offered at time 0 for detecting the failure at the n th step of s ;
- $t[n, s]$: time elapsed till the moment when module $s[n]$ is examined at step n of sequence s ;
- $P[n, s]$: the probability of detecting the failure at the n th step of sequence s (if no discount is involved).

Then the time $T[n, s]$ spent to detect and localize the failure at the n th step of s will be, as before:

$$T[n, s] = \sum_{d=1, \dots, n} t[d, s].$$

The reward discounted before the failure is localized will be:

$$W[n, s] \exp(-dT[n, s]),$$

while the improved probability of detecting the failure at the n th step of s will be:

$$P[n, s] = \prod_{i=1, \dots, N} (b_i)^{u(i, n, s)}.$$

Then the total expected discounted reward when using the sequence s will be:

$$\text{Rew}(s) = \sum_{n=1, \dots, \text{inf}} P[n, s] \prod_{i=1, \dots, N} (b_i)^{u(i, n, s)} D(n, s) \exp(-dT[n, s]).$$

The following generalization of the Sweat Rule [12] takes place:

Theorem 3. *Let the values of ratio is*

$$R_{ij} = b_i W_i Q_{ij} \exp(-dt_i) / [1 - b_i \exp((-dt_i))]$$

for all i and j be arranged in non-increasing order, and let s^ be such that if the ratio R_{ij} is the n th largest one in the ordering*

then the n th step of s^* is the j th search of module i . Then s^* is optimal.

Optimality conditions for Min-Loss Search

Consider a more general situation when the search times depends on the “history” of the search: $t_i = t_i(n, s)$. Let f be an arbitrary loss functions of the total search time

$$T[n, s] = \sum_{d=1, \dots, n} t[d, s]$$

at step n , and $F(s)$ be the corresponding total expected losses:

$$F(s) = \sum_{d=1, \dots, \inf} P[n, s] c[n, s] f(T[n, s]).$$

Definition 1. Let $IF = IF_i(k)$ be a real-valued function depending on the index i , the search step k and the current input data characterizing the module i ; this function is called the *index function*, or the *fitness function*. A local search policy is called an index policy generated by $IF = IF_i(k)$ if at each search step one inspects a module with currently greatest value I_i , among all i .

Definition 2. An objective function $F(S)$ permits an index policy if there exists a single-index fitness function $IF_i = IF_i(k)$ such that for all values of N and any input data the index policy generated by $IF_i = IF_i(k)$, is optimal.

Proposition 1. Let $f_i(t)$ be strictly increasing and sufficiently smooth (there exists the 3rd derivative) function on R_+ . The expected loss function $F(s)$ permits an index policy if and only if (i) $f_i(t) = a_i(t) + b_i$, or (ii) $f_i(t) = a_i \exp(kt) + b_i$.

Proposition 2. Let $f_i(t)$ be strictly increasing and sufficiently smooth (there exists the 3rd derivative) function on R_+ . The expected reward function $R(s)$ permits an index policy if and only if (i) $f_i(t) = a_i t + b_i$, or (ii) $f_i(t) = a_i \log(kt) + b_i$, for $i = 1, \dots, N$.

If we compare the properties of these search algorithms with the conceptual model of fault diagnostic described in the Section 2, we can see that classes $(PP|s|pr|imperf|Be, \text{Min-cost})$ and $(PP|s|pr|imperf|M, \text{Max-reward})$ cover very wide part of the area of problems concerning the complex computer systems maintenance [1, 8, 16, 17].

4. Conclusion and future work

The considered classification allows us to define various classes of search models, which differ in their assumption about character of search and possible failures during search (e. g., whether an algorithm is “noisy-free” or “noisy” from the point of view of the objects goodness identification). A consideration of relationship between the classes can show how we should choose (or change) a search model for a given conditions. It may be useful to manage the process of diagnostic planning. For example, we can see that the properties of search algorithms (attributes **FM**, **FOM**, Section 2) affect the parameters of fault search model rather than other fault characteristics. Further, the relationships allow us, for example, establish a hierarchy of both search conceptual models and the search algorithms, where the hierarchy levels may correspond to the degree of generality of the models (that is the more complex diagnostic figure (fault search conditions, testing scenario, etc. the higher level of the search model), or, the more aspects of the search strategy are taken into account in a search model, the higher its level).

Such characterization may be very useful in teaching for large computer system (including networks) maintenance and service, since allow to co-ordinate in a clear manner technical and mathematical issues. Besides, such classification could be used in the field of fault localization and diagnosis by expert systems [14], where rule-based representations of their knowledge base are in use. Using the above approach, it is possible to perform experiments and solve practical classes of search problems with uncertain data and under different working scenarios.

The following unresolved currently issues are very topical for this research area:

1. To extend the model: when costs depend in time, the learning decreases the overlooking probability;
2. When **K** contains several failed modules;
3. When search at each step considers several modules (and their number may vary in time).

References

1. Li H., Baras J. S. and Mykoniatis G. An Automated, Distributed, Intelligent Fault Management System for Communication Networks // ATIRP'99, 2–4 February, 1999.

2. *Steinder M., Sethi A.S.* A survey of fault localization techniques in computer networks // *Science of Computer Programming*. 2004. V. 53. P. 165–194.
3. *Alonso L., Chassaing P., Reingold E.M., and Schott R.*, The Worst-Case Chip Problem // *Info. Proc. Let.* 2004. V. 89. P. 303–308.
4. *Brodie M. et al.* Intelligent probing: A cost-effective approach to fault diagnosis in computer networks // *IBM System Journal*. 2002. V. 41, No. 3.
5. *Pearl J.* Probabilistic reasoning in intelligent systems: networks of plausible inference. — Palo Alto, Calif., U.S.A.: Morgan Kaufmann, 1988.
6. *Pablo H., Ibarguengoytia L.* Enrique Sucar & Eduardo Morales A Probabilistic Model Approach for Fault Diagnosis // Eleventh International Workshop on Principles of Diagnosis, Morelia, Mexico, June, 2000. P. 79–86.
7. *Shayman M.A.* Risk-Sensitive Decision-Theoretic Diagnosis // *IEEE Transactions on Automatic Control*. July 2001. V. 46. P. 1166–1171.
8. *Steinder M., and Sethi A.S.* Probabilistic Fault Localization in Communication Systems Using Belief Networks // *IEEE/ACM Transactions on Networking*. October 2004. V. 12, No. 5.
9. *Schulz M.H., Trischler E., and Sarfert T.M.* Socrates: A highly efficient automatic test pattern generation system // *IEEE Transactions on circuits CAD*. January 1988. P. 126–137.
10. *Siewiorek D.P., Maxon R.A., and Narasimhan P.* Experimental research in dependable computing: From Faults to Manifestations // Technical report of CS department of CMU, www.cs.cmu.edu/~faculty/ThinkPieces/Siewiorek-Maxon-Narasimhan.doc.
11. *Ahlsvede R., Weneger I.* Search Problems. — N.Y., Wiley, 1987.
12. *Levner E.* Infinite-horizon scheduling algorithms for optimal search for hidden objects // *International transactions in operational research*. 1994. V. 1, No. 2. P. 241–250.
13. *Wegener I. mit Lösner U.*: Discrete sequential search with positive switch cost // *Mathematics of Operations Research*. 1982. Bd. 7. S. 426–440.
14. *Patel A., McDermott G., Mulvihill C.*, Integrating network management and artificial intelligence // In: B. Meandzija, J. Westcott (Eds.), *Integrated Network Management*, North-Holland, Amsterdam, 1989. P. 647–660.
15. *Amari S.V., Pham H., and Dil G.*, Optimal Design of k-out-of-n: G Subsystems Subjected to Imperfect Fault-Coverage // *IEEE Transactions on Reliability*. December 2004. V. 53, No. 4. P. 567.

16. *Alcaide D., Levner E., Frenkel S., and Zakharov V.* On-line search algorithms for fault diagnostics in large-scale computer communication networks // Proc. of the XXIX Spanish Meeting of Statistics and Operations Research, SEIO-2006, La Laguna, Spain, 2006.
17. *Frenkel S., Levner E., Zakharov V.* Characterization of probabilistic faults diagnostic models // Proc. IEEE Information Technology, Research and Education 2006, ITRE'06, International Conf., 2006. P. 156-160.
18. *Kendall D.G.* Some problems in the theory of queues // Journal of the Royal Statistical Society, Series B, 1951. P. 151–185.

COMPLEXITY AND CONSISTENCY OF STATISTICAL CRITERIA

Alexander Grusho

Lomonosov Moscow State University, Moscow, Russia

Nick Grusho

Russian State University for the Humanities, Moscow,
Russia

Elena Timonina

Russian State University for the Humanities, Moscow,
Russia

We investigate relations between consistency of statistical criteria sequences in finite probability spaces and asymptotic complexity of test calculations. It is proved that for every consistent test sequence it is possible to construct another consistent sequence for which complexity of calculations of membership functions of new criteria critical sets can be made asymptotically small in comparison with similar complexity for initial criteria. However such simplification appears to be fictitious as a matter of fact. To prevent fictitious simplification of calculations of membership functions of critical sets in a sequence of criteria it is necessary to set additional restrictions on classes of considered criteria. It is shown, that in case of such natural restrictions any simplification of calculations can lead to a failure of a consistency of criteria sequence.

Estimations of efficiency of two-level criteria are resulted, when simple but not consistent in the given class of alternatives criteria work at first, and only in case of nonacceptance of a hypothesis hard computable but consistent criteria are applied.

Let $X = \{x_1, \dots, x_m\}$ be a finite set, X^∞ be a set of infinite sequences, in which elements belong to X , \mathcal{A} be σ -algebra, generated by cylindrical sets, P_0 be a probability measure on the measurable space (X^∞, \mathcal{A}) . Furthermore, there is a family of probability measures P_ϑ , $\vartheta \in \Theta$, on the same measurable space.

This work was supported by the Russian Foundation for Basic Research (grants 07-01-00484 and 07-07-00236).

This work was presented at III International Workshop "Applied Problems of Probability Theory and Mathematical Statistics related to modeling of information systems" (Aosta, Italy, January 2008).

Consider projections of the introduced probability measures onto the first n coordinates of sequences from X^∞ and denote them accordingly $P_{0,n}$ and $P_{\vartheta,n}$. Concerning the given measures for every n we define a problem of testing of a statistical hypothesis $H_{0,n} : P_{0,n}$ against complex alternative $H_{1,n} : \{P_{\vartheta,n}, \vartheta \in \Theta\}$. A criterion of significance value α is described by a critical set S_n , $P_{0,n}(S_n) \leq \alpha$ and a power of criterion $W_n(\vartheta) = P_{\vartheta,n}(S_n)$. The sequence of statistical criteria with critical sets S_n is named consistent [4], if for all $\alpha > 0$ the power of criteria $W_n(\vartheta) \rightarrow 1$ for every $\vartheta \in \Theta$.

A statistical criterion can be described by a membership function of the critical set of a criterion. In discrete mathematics we can speak about complexity of a membership function of the given set. Thus, it is possible to speak about complexity of calculation of a statistical criterion. In the considered above model we have a sequence of statistical criteria for every n and $n \rightarrow \infty$. That is why we say about asymptotic estimations of complexity of a criteria sequence.

A relation between asymptotic statistical properties of a criteria sequence and asymptotic complexity of their calculations was rarely investigated in the discrete mathematics. It is natural to consider properties of consistency of a criteria sequence and asymptotic complexity of this sequence of criteria. For example, in the work [2] it was shown, that the problem of detection of the chosen set of nodes in a random graph had a consistent decision when power of the given set of nodes grows as $C \ln n$, where n is the number of nodes of the random graph. However the consistent sequence of criteria has asymptotically high complexity. If we consider simply computable criteria the corresponding border of detectability of the chosen set becomes equal $C\sqrt{n}$.

We investigate the next questions. Can any reduction in complexity lead to loss of a consistency? Can we get reduction in complexity of a sequence of criteria and preservation of consistency?

We define a one-dot distribution as a singular distribution of probabilities concentrated in one point of space X^∞ . Let P_0 be a one-dot distribution concentrated in the point $\omega = (x_1, x_1, \dots, x_1, \dots)$, all coordinates of which are equal to x_1 . As alternatives we consider a set of one-dot distributions $\{P_{\vartheta}, \vartheta \in X^\infty, \vartheta \neq \omega\}$. As the basic operation for an estimation of complexity we shall take comparison of the next sign of observed sequence with x_1 . For every n the projection $P_{0,n}$ is a singular

distribution on X^n , concentrated in the point ω_n , where ω_n is a vector of length n with all coordinates are equal to x_1 . Each of alternatives $P_{\vartheta,n}$ is also a singular distribution, concentrated in the point ϑ_n . As the critical set for testing of the hypothesis $H_{0,n}$ against complex alternatives $H_{1,n}$ we put $S_n = X^n \setminus \{\omega_n\}$. For this criterion $P_{0,n}(S_n) = 0$ and $P_{\vartheta,n}(S_n) = 1$ for any ϑ , for which projection $\vartheta_n \neq \omega_n$. Complexity of calculation of the given criterion does not exceed n for all points of space X^n .

Thus defined sequence of criteria is consistent. Really

$$\lim_{n \rightarrow \infty} P_{0,n}(S_n) = 0,$$

and for everyone $\vartheta \in X^\infty$, $\vartheta \neq \omega$ there exists n such, that $\vartheta_n \neq \omega_n$. Then, since this n , the power $W_n(\vartheta) = 1$.

Complexity of calculation of any other criterion can be lowered only because we shall not compare each coordinate of observed vector from X^n with unique admissible value x_1 . Then if complexity of a sequence of statistical criteria is equal $o(n)$, it means that the membership function of the critical set does not depend on the significant part of coordinates. Let's assume that the sequence of sets of skipped coordinates is monotone nondecreasing with growth n . Then the sequence of critical sets D_n for such sequence of criteria will possess following property

$$\lim_{n \rightarrow \infty} [(X^n \setminus D_n) \times X^\infty] \supset \{\omega\}.$$

That is for all n cylindrical sets, generated by critical sets of criteria, will not possess at least one more sequence $\vartheta \neq \omega$. It means that for this sequence ϑ power of the test does not tend to 1.

Thus we have proved, that reduction in complexity of criterion leads to loss of property of consistency for corresponding sequence of criteria.

We'll show that for an arbitrary consistent sequence of criteria there is another consistent sequence of criteria against the same alternatives with much more lower complexity.

Theorem 1. *For an arbitrary consistent sequence of criteria for testing of $H_{0,n}$ against $H_{1,n}$ with critical sets S_n , $n = 1, 2, \dots$, and complexity of their membership functions $g(S_n)$, where $g(S_n) \rightarrow \infty$ if $n \rightarrow \infty$, there exists a consistent sequence of criteria for testing of $H_{0,n}$ against $H_{1,n}$, complexity of which is asymptotically small in comparison with complexity of criteria of reference.*

Proof. By the condition of consistency $\forall \alpha \in (0; 1]$

$$P_{0,n}(S_n) \leq \alpha$$

and for every $\vartheta \in \Theta$

$$\lim_{n \rightarrow \infty} P_{\vartheta,n}(S_n) = 1.$$

Then

$$P_0(S_n \times X^\infty) \leq \alpha$$

and for any $r = 1, 2, \dots$,

$$P_0(S_n \times X^r \times X^\infty) \leq \alpha, \quad P_{0,n+r}(S_n \times X^r) \leq \alpha.$$

Let $k = k(n) = o(n)$, $k(n) \rightarrow \infty$ if $n \rightarrow \infty$ and $g(S_{k(n)}) = o(g(S_n))$. Let's consider criteria with a sequence of critical sets D_n , $n = 1, 2, \dots$. Here for $k(n) \geq n$

$$D_n = S_{k(n)},$$

and for $k(n) < n$

$$D_n = S_{k(n)} \times X^{n-k(n)}.$$

The constructed sequence of criteria is consistent, because

$$\begin{aligned} P_{0,n}(D_n) &= P_0(D_n \times X^\infty) = P_0((S_{k(n)} \times X^{n-k(n)}) \times X^\infty) \\ &= P_{0,k(n)}(S_{k(n)}) \leq \alpha, \end{aligned}$$

$$\begin{aligned} P_{\vartheta,n}(D_n) &= P_{\vartheta}(D_n \times X^\infty) = P_{\vartheta}((S_{k(n)} \times X^{n-k(n)}) \times X^\infty) \\ &= P_{\vartheta,k(n)}(S_{k(n)}) \xrightarrow{n \rightarrow \infty} 1. \end{aligned}$$

Complexities of these criteria are equal to $g(S_{k(n)})$ because the critical sets do not depend on $n - k(n)$ last coordinates. The theorem is proved.

The considered above example and the proved theorem seem to be contradictory. In the example we demanded monotony of sequence of sets of the unseen coordinates, and in the theorem there is no this condition. We shall name this contradiction as the conflict of "laziness". Really, by consideration of infinite sequence of decisions the observer can perform the work on a step n as though he is on a step $k < n$ and he does not look through $n - k$ last coordinates of observed vector. In view of a consistency he knows, that sooner or later he will make all the work with the

same asymptotic result. Complexity connected with “laziness” does not reflect real complexity of the sequence of criteria.

To avoid the conflict it is necessary to demand additional conditions for sequences of criteria. We shall present two examples of such conditions.

Each critical set S_n can be constructed from the set S_{n-1} with the help of withdrawal of some sequences from the set $S_{n-1} \times X$. Then the complexity of the criterion $g(S_n)$ which is constructed with the help of the specified algorithm satisfies to the difference equation

$$g(S_n) = g(S_{n-1}) + f(n),$$

where $f(n)$ is the complexity of withdrawal of sequences from the set $S_{n-1} \times X$.

Monotone nonincreasing sequence of the closed cylindrical sets $S_n \times X^\infty \supseteq S_{n+1} \times X^\infty$, $n = 1, 2, \dots$ corresponds to the sequence of critical sets S_n , $n = 1, 2, \dots$. Then there is a limit

$$S = \bigcap_{n=1}^{\infty} (S_n \times X^\infty)$$

and the set

$$A = X^\infty \setminus S$$

is the open set in Tychonoff product [1, 5].

For the sequence of critical sets S_n , $n = 1, 2, \dots$, it is possible to consider, that

$$\lim_{n \rightarrow \infty} P_{0,n}(S_n) = 0.$$

From here it follows, that $P_0(S) = 0$ and $P_0(A) = 1$.

Let's select a class of alternatives $\Theta_1 \subseteq \Theta$ for which there are sets $A(\vartheta)$, $\vartheta \in \Theta_1$, such that $A(\vartheta) \cap A = \emptyset$ and $P_\vartheta(A(\vartheta)) = 1$. For this class of alternatives the sequence of criteria with critical sets S_n , $n = 1, 2, \dots$, is consistent [3].

Calculation of membership function of the critical set S_n can be replaced with the membership function of the set complement $X^n \setminus S_n$. That is if the observed value does not belong to the set complement S_n then the observed value belongs to S_n . In some cases such calculation of membership functions is easier.

Let's denote $D_n = X^n \setminus S_n$, $n = 1, 2, \dots$. Let D_n can be constructed from the set $D_{n-1} \times X$ with the help of withdrawal of

some sequences and let $g(D_n)$ be a complexity of realization of this algorithm. Then for some function $f(n)$ we have difference equation

$$g(D_n) = g(D_{n-1}) + f(n),$$

where $f(n)$ characterizes complexity of withdrawal of sequences from the set $D_{n-1} \times X$. In this case we have monotone nonincreasing sequence of cylindrical sets

$$D_n \times X^\infty \supseteq D_{n+1} \times X^\infty, \quad n = 1, 2, \dots$$

For the sequence of critical sets S_n , $n = 1, 2, \dots$, it is possible to consider, that

$$\lim_{n \rightarrow \infty} P_{0,n}(S_n) = 0.$$

From here it follows, that $P_0(D) = 1$.

The set D is closed in Tychonoff product [1, 5].

We shall select a class of alternatives $\Theta_1 \subseteq \Theta$ for which there are sets $A(\vartheta)$, $\vartheta \in \Theta_1$, such that $A(\vartheta) \cap D = \emptyset$ and $P_\vartheta(A(\vartheta)) = 1$. Against this class of alternatives the sequence of criteria with critical sets S_n , $n = 1, 2, \dots$, is consistent [3].

When the critical sets or their sets complement are obtained as the Cartesian products $S_n = C^n$ of some subset $C \subseteq X$, then $f(n) = \text{const}$. In this case complexity $g(S_n)$ is linear function of n .

At the analysis of complexity the following simple theorem is required to us.

Theorem 2. *Let us consider some classes of alternatives $H_{1,n}^{(k)} : P_{\vartheta,n}$, $\vartheta \in \Theta_k$, $k = 1, 2, \dots, M$, and $H_{1,n} : P_{\vartheta,n}$, $\vartheta \in \bigcup_{k=1}^M \Theta_k$. There are consistent sequences of criteria for testing of $H_{0,n}$ against each of alternatives $H_{1,n}^{(k)}$ iff there is a consistent sequence of criteria for testing of $H_{0,n}$ against $H_{1,n}$.*

Proof. Existence of a consistent sequence of criteria for testing of $H_{0,n}$ against each alternative $H_{1,n}^{(k)}$, $k = 1, 2, \dots, M$, means, that there is a sequence of sets $S_n^{(k)}$ such, that $P_{0,n}(S_n^{(k)}) \rightarrow 0$ and for every $\vartheta \in \Theta_k$ it is carried out $P_{\vartheta,n}(S_n^{(k)}) \rightarrow 1$. We shall denote $S_n = \bigcup_{k=1}^M S_n^{(k)}$. Then

$$P_{0,n}(S_n) \leq \sum_{k=1}^M P_{0,n}(S_n^{(k)}).$$

Thus, the sequence of criteria with critical sets S_n satisfy the condition

$$P_{0,n}(S_n) \longrightarrow_{n \rightarrow \infty} 0.$$

For $\vartheta \in \Theta_k$ it is carried out

$$P_{\vartheta,n}(S_n) \geq P_{\vartheta,n}(S_n^{(k)}) \longrightarrow_{n \rightarrow \infty} 1.$$

Let S_n be a sequence of critical sets of consistent sequence of criteria for testing of $H_{0,n}$ against alternative $H_{1,n}$. The same sequence of criteria is consistent for testing of $H_{0,n}$ against $H_{1,n}^{(k)}$, $k = 1, 2, \dots, M$. The theorem is proved.

Let's consider a special case of application of the theorem 2. Assume that the sequence of criteria of a significance value $\frac{\alpha}{2}$ with critical sets $S_n^{(1)}$ has asymptotic complexity $g(S_n^{(1)})$, and the sequence of criteria of a significance value $\frac{\alpha}{2}$ with critical sets $S_n^{(2)}$ has asymptotic complexity $g(S_n^{(2)})$.

Assume, that $g(S_n^{(1)}) \ll g(S_n^{(2)})$ when $n \rightarrow \infty$. At the same time the first sequence of criteria is consistent only for a set of alternatives $\Theta_1 \subset \Theta$, and the second sequence of criteria is consistent for the whole set Θ .

For reduction of complexity we advise the next procedure of testing of $H_{0,n} : P_{0,n}$ against alternatives $H_{1,n} : P_{\vartheta,n}$, $\vartheta \in \Theta$.

First we test a hypothesis $H_{0,n}$ with the help of criterion $S_n^{(1)}$, and in case of its deviation this hypothesis is checked with the help of criterion $S_n^{(2)}$. There significance value is α .

Complexity of the constructed procedure is equal to

$$g = g(S_n^{(1)}) + g(S_n^{(2)}) \cdot I(S_n^{(1)}),$$

where $I(A)$ is the indicator of event A .

Expectation of the complexity of the constructed procedure in the measure $P_{0,n}$ is estimated in the following way

$$E_0 g \leq g(S_n^{(1)}) + \frac{\alpha}{2} \cdot g(S_n^{(2)}).$$

This formula shows, what the prize we can get, applying simple criterion at first, and then specifying received decision with the help of more complex criterion.

References

1. *Bourbaki N.* Topologie Générale. — Russian translation, Moscow, Science, 1968. P. 272.
2. *Grusho A. A.* Some statistical graph problems // *Mathematical notes*. 1984. V. 36, No. 2.
3. *Grusho A. A., Timonina E. E.* Some relations between discrete statistical problems and properties of probability measures on topological spaces // *Discrete Mathematics and Applications*. 2006. V. 16, No. 6. P. 547–554.
4. *Lehmann E. L.* Testing Statistical Hypotheses (Springer Texts in Statistics). — Springer, 2nd edition, 1997. P. 600.
5. *Neveu J.* Bases mathématiques du calcul des probabilités. — Paris, Masson, 1964. P. 310.

BAYESIAN QUEUEING AND RELIABILITY MODELS¹

Alexey Kudryavtsev

Faculty of Computational Mathematics and Cybernetics,
Lomonosov Moscow State University, Moscow, Russia

Sergey Shorgin

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

Vsevolod Shorgin

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

Valery Chentsov

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

The Bayesian approach for certain tasks of queueing systems theory and reliability theory is investigated. The method provides the randomization of system characteristics with regard of a priori distributions of input parameters. This approach could be used, for instance, for calculating average values and for construction of confidential intervals applicable for performance and reliability characteristics of large groups of systems or devices. The results for certain models of input flow and service time parameters are presented.

1. Introduction and main assumptions

Theory of queueing systems is a well-developed mathematical discipline. Based on it a substantial number of positive R&D results have been generated. The results obtained in studying queueing systems and networks proved to be of significant profundity and importance from mathematical and practical points of view. In fact queueing systems and networks are able to model a broad class of real systems, info-telecommunication systems and networks being

¹ This work was supported by the Russian Foundation for Basic Research (grant 08-07-00152).

This work was presented at III International Workshop “Applied Problems of Probability Theory and Mathematical Statistics related to modeling of information systems” (Aosta, Italy, January 2008).

in the first place. In order to reflect real processes in a more adequate way, the present development of queueing theory is being carried out mostly with a focus on studying more complex service disciplines, input flows and service time distributions with more and more complicated probabilistic characteristics.

One of the directions of generalization of problem formulations is the complication of probabilistic structure of one or more queueing systems input parameters. Instead of considering traditional input flows, the researchers study Cox flows, self-similar flows, Markovian and semi-Markovian flows, etc. Similar generalizations are made regarding service times distributions. To some extent, these generalizations can be interpreted as the randomization result of these or those parameters of more “simple” flows and service times distributions. Thus, Cox process is obtained as a result of special randomization of Poisson flow intensity, etc.

All these generalized modern formulations assume that stochastic method of randomization “affects” the parameters of a system precisely during its functioning, meaning that we primarily know the kind of the system we are “dealing with”, even when the system is rather complicated and then we study characteristics of this particular “primarily fixed” system. However, in real life often the system under study is specified in some sense vaguely, or inaccurately. For example, even when we deal with the simplest systems of $M|G|1$ type, we may not know a priori the input flow parameter λ and the service parameters μ and σ^2 . Such situations can occur studying the whole class of queueing systems when the only known characteristics are the input flow types the service distribution and the service discipline, but at the same time the concrete parameters of these flows and distributions, generally speaking, vary for different queueing systems of a given class. A researcher does not know a priori the queueing system belonging to the given class he is dealing with. For example, such situation can take place testing a series of uniformed commutation or transmission devices manufactured by the same company. Spread in some of their performances can be explained by natural technological deviations during manufacturing process. In this particular case, since the unknown characteristics are the “initial” parameters of the flows and service times, a natural thing could be the use of a randomized approach according to which the values λ , μ and σ^2 become the elements of a probabilistic space, but in general, one can

speak about probabilistic space with uniformed queueing systems being its elements. In this situation it is quite natural that the calculated characteristics of such randomized queueing system are randomization of similar characteristics of “usual” queueing system of similar type taking into account a priori distribution of queueing system input parameters.

So, in the same example concerning a $M|G|1$ queueing system there arise the tasks of “common” characteristics randomization of such systems with regard for a priori input parameters distributions. In other words, we can make assumption about exponential, uniform or any other distribution of one or several values λ , μ or σ^2 (that become random variables under such approach), about their dependence or independence, etc. The obtained results could be used, for instance, to calculate “in general” average values and to construct confidential intervals applicable for these or those characteristics of the queueing system class under consideration. Naturally, such approach queueing models development can be called Bayesian and it was formulated for the first time in [1].

The Bayesian approach can be used also in problems of reliability estimation. As it is known (see [2]), the availability factor of the restorable device in a stationary mode can be calculated using the formula

$$k = \frac{\lambda^{-1}}{\lambda^{-1} + \mu^{-1}} = \frac{\mu}{\lambda + \mu},$$

where λ^{-1} is the average operating time between failures, and μ^{-1} is the average restoration time. If we accept the hypothesis stated above that the device under consideration is randomly selected from some set of similar devices whose average reliability characteristics vary, then according to the reasonings presented above, values λ and μ could be considered as random. Hence, under these assumption the availability factor k is random, too, and its distribution depends on distributions of values λ , μ . The obtained results in this field could be used, for instance, for calculating “in general” average values and for the construction of confidential intervals for reliability characteristics of the overall set of investigated devices.

2. The Bayesian approach to queueing systems

In order to explain the essence of the task formulation we present the following example. Let us consider a situation when an observer deals with rather large series of queueing systems

$M|M|1|0$ that differ only in service distribution parameter. In particular, these can be certain devices, commutators, routers or any other servicing tools. It is known in advance that their functioning can be modelled by a system belonging to the above-described type., i.e these systems have identical service discipline, types of input flow and of service times distribution.

This example assumes that the input flow characteristics are also identical for all the systems of a given series; only numerical characteristics of *service* are different (i. e. the parameters of exponential distribution).

Dispersion in characteristics of service is due to technological (design) reasons and the main aspect of the problem statement is the fact that the researcher does not know what the real value of service parameter of the system belonging to a given series under study that was selected by him at random. The only thing that he knows is “a priori” distribution of this parameter (since the series is supposed to be large, one can consider stochastic phenomena in relation with that series and introduce probabilistic distributions). The researcher is interested in finding out service characteristics for a series as a whole (or characteristics of the system “selected at random”). Obviously, along with traditional factors of stochasticity that occur in queueing systems (stochasticity of input flow and service processes), there appears one more factor of stochasticity related to *randomized selection of the system under study*.

Let us assume that the service parameter μ of the systems under study can take only two values: μ_1 and μ_2 with probability p_1 and p_2 , respectively. In “physical terms” it means that among the system series under study (routers, machine tools, etc.) only two “varieties” of servicing devices occur. Devices belonging to the first variety provide the service with parameters μ_1 , while devices of the second variety provide the service with parameter μ_2 . Then the loading factor of the system “selected at random” becomes the random variable that takes the values λ/μ_1 with probability p_1 and λ/μ_2 with probability p_2 . The steady-state probability of blocking the “selected” system due to the interference of the random factor of selecting a concrete system becomes “random” itself and takes the values $\lambda/(\lambda + \mu_1)$ with probability p_1 (it is the probability that a system belonging to the first variety has “fallen into the researcher’s hands”) and $\lambda/(\lambda + \mu_2)$ with probability p_2 (meaning that a system of the second variety “has fallen into the researcher’s hands”). It is

natural that the “averaged” blocking probability of such “Bayesian” queueing system is equal to $p_1\lambda/(\lambda + \mu_1) + p_2\lambda/(\lambda + \mu_2)$.

As we can see, there is no need to use the methods of queueing theory for studying the Bayesian queueing systems. Bayesian system is “randomization” of a certain “ordinary” queueing system, meaning that the Bayesian queueing system characteristics can be calculated by means of randomizing subsequent averaging (by a priori distribution of the parameter or parameters) of the “ordinary” queueing system characteristics that have been calculated earlier by using the methods of queueing theory. In other words, the mathematical part of the job comes to this particular randomization and averaging. At the same time, it is an expedient from both technological and conceptual points of view to accomplish randomization of stationary characteristics of “ordinary” queueing systems and obtain the steady-state characteristics of Bayesian queueing systems.

We would like to point out one more substantial model that can be described mathematically with the help of Bayesian queueing system. Let’s assume that a researcher considers not a series of systems with quantitative parameters that change with the time. For example, there exists a servicing device, one of its elements being replaced by another one at the moments that we do not know beforehand, then being replaced by the third one, etc. Such a system can be the frontier post at the airport, where an officer on duty is relieved from time to time at the moments not known by the observers (passengers). The only things an observer knows are the probability that he will have “come upon” a certain concrete frontier officer and an average time of passport checking by each frontier.

Under such approach the system structure and service discipline do not change with the time while only quantitative parameter of distribution of service changes (e.g. intensity). The input flow parameter can change in a similar way. There is no information about the moment when changes occur. The researcher is aware only of *distribution* of the values of “changeable”, random parameters he “comes across” while examining the system at a “random” moment of time.

Since it is assumed that the researcher does not have any information about the moments of the system “reorganization”, and even about distribution of these moments, it is impossible to describe transient processes within such kind of a system. Therefore, it is possible to carry out analysis (and subsequent randomization)

of only *steady-state* distributions of the queueing system under analysis. In order to give meaning to this problem statement, it is necessary to make an assumption that the system changes quite “rarely” so that at each interval of constancy of the parameters, the queueing system “had time” to reach steady-state condition. Of course, the results of such analysis will be rough because steady-state condition, strictly speaking, cannot be reached in real life.

3. Simple models of “Bayesian” queueing systems

Below two more simplest models of “Bayesian” queueing systems are presented in order to provide further elucidation of specific character of the problems that emerge under such an approach and of the obtained results. The results of this chapter were presented in [3]

3.1. Uniform distribution of λ and μ : loading factor. Let us consider an arbitrary queueing system with input flow intensity λ and service intensity μ . The loading of such system is equal to $\rho = \lambda/\mu$. As it is generally known, the availability of steady-state mode of the system under consideration depends on the value ρ which apperas in many formulae that describe characteristics of different queueing systems. Hence, the study of the value ρ should be considered within the frames Bayesian theory of queueing systems.

The variety of possible and interesting distributions of variables λ and μ for their joint applications is rather wide. We consider one of the simplest but at the same time very common in practice cases when the values λ and μ are independent and uniformly distributed on some certain pre-determined segments. Such model is good for describing situations when some legitimate interval of values have been assigned for both values λ and μ (or for any of them), but the real value λ or/and μ can vary within such limits.

Assume that the random variable λ has a uniform distribution on the segment $[a_\lambda, b_\lambda]$, the random variable μ has a uniform distribution on $[a_\mu, b_\mu]$, with $0 \leq a_\lambda \leq b_\lambda$, $0 \leq a_\mu \leq b_\mu$.

In this case, the cumulative function of the random variable $\rho = \lambda/\mu$ distribution can be written down as follows:

$$P\{\rho < x\} = \iint_{\lambda/\mu < x} \frac{1}{b_\lambda - a_\lambda} \frac{1}{b_\mu - a_\mu} d\lambda d\mu.$$

Subsequent calculations depend essentially on relation between the values a_λ/a_μ and b_λ/b_μ . Let us suppose for the sake of definiteness that $a_\lambda/a_\mu \leq b_\lambda/b_\mu$. Then:

provided $x \leq a_\lambda/b_\mu$

$$P\{\rho < x\} = 0,$$

provided $a_\lambda/b_\mu \leq x \leq a_\lambda/a_\mu$

$$P\{\rho < x\} = K \frac{(b_\mu x - a_\lambda)^2}{2x},$$

provided $a_\lambda/a_\mu \leq x \leq b_\lambda/b_\mu$

$$P\{\rho < x\} = K \left(\frac{a_\mu + b_\mu}{2} x - a_\lambda \right) (b_\mu - a_\mu),$$

provided $b_\lambda/b_\mu \leq x \leq b_\lambda/a_\mu$

$$P\{\rho < x\} = 1 - K \frac{(b_\lambda - a_\mu x)^2}{2x},$$

provided $x \geq b_\lambda/a_\mu$

$$P\{\rho < x\} = 1,$$

when

$$K = \frac{1}{(b_\mu - a_\mu)(b_\lambda - a_\lambda)}.$$

Let us derive the density of random variable ρ :

provided $x \leq a_\lambda/b_\mu$

$$f_\rho(x) = 0,$$

provided $a_\lambda/b_\mu \leq x \leq a_\lambda/a_\mu$

$$f_\rho(x) = K \left(\frac{b_\mu^2}{2} - \frac{a_\lambda^2}{2x^2} \right),$$

provided $a_\lambda/a_\mu \leq x \leq b_\lambda/b_\mu$

$$f_\rho(x) = K \left(\frac{b_\mu^2 - a_\mu^2}{2} \right),$$

provided $b_\lambda/b_\mu \leq x \leq b_\lambda/a_\mu$

$$f_\rho(x) = K \left(\frac{b_\lambda^2}{2x^2} - \frac{a_\mu^2}{2} \right),$$

provided $x \geq b_\lambda/a_\mu$

$$f_\rho(x) = 0.$$

Through accomplished elementary calculations, we derive the average value and the second moment of random variable ρ , that are respectively equal to:

$$\mathbf{E}\rho = \frac{b_\lambda + a_\lambda}{2(b_\mu - a_\mu)} \ln \frac{b_\mu}{a_\mu},$$

$$\mathbf{E}\rho^2 = \frac{a_\lambda^2 + a_\lambda b_\lambda + b_\lambda^2}{3a_\mu b_\mu}.$$

It is evident that if $b_\lambda - a_\lambda \rightarrow 0$ and $b_\mu - a_\mu \rightarrow 0$, i. e. contracting the range of the random variable λ to some fixed point λ_0 , and the range of the random variable μ to some fixed point μ_0 , the value $\mathbf{E}\rho$, as it should be, tends to λ_0/μ_0 , and the value $\mathbf{E}\rho^2$ tends to λ_0^2/μ_0^2 .

Moreover, we note that the dependence of the average value of ρ on distribution λ is reduced to dependence on the mathematical expectation λ . At the same time, dependence of $\mathbf{E}\rho$ on parameters of distribution μ has a more complex look.

In the case $a_\lambda/a_\mu \geq b_\lambda/b_\mu$, the formulae for calculating the cumulative and density functions of the random variable ρ are similar. The mathematical expectation and the second moment of the random variable ρ in this particular case coincide with the values that have been calculated previously.

Based on the obtained results, it would be easy to calculate other necessary characteristics of value ρ .

It is worthwhile to observe that the examined model allows to study an important situation when $\lambda < \mu$ has the probability 1. In this case $\rho < 1$, which is the condition of ergodicity of the systems having one servicing device. By virtue of postulated independence of random values λ and μ , and the condition for $\lambda < \mu$ is satisfied only if the condition $0 \leq a_\lambda \leq b_\lambda \leq a_\mu \leq b_\mu$ holds.

3.2. Exponential λ and μ distribution: loading factor, probability of losses in the system $\mathbf{M|M|1|0}$ and availability factor.

Let us consider another probabilistic model for the values λ and μ . In a situation when there is no a priori information about their mean values, it we can consider as a “first approximation” a model where λ and μ are exponentially distributed with known averages, $1/l$ and $1/m$ respectively). Assumption about λ and μ has been retained.

So, the cumulative function of the random variable λ distribution is equal to $1 - \exp(-lu)$ and the cumulative function of

the random variable μ distribution is equal to $1 - \exp(-mu)$. As we did in the previous section, let us first of all consider $\rho = \lambda/\mu$. Obviously, for $x \geq 0$ we get

$$\begin{aligned} P\{\rho < x\} &= P\{\lambda < \mu x\} = \int_0^\infty P\{\lambda < xy\} dP\{\mu < y\} \\ &= \int_0^\infty [1 - \exp(-lxy)] m \exp(-my) dy = \frac{lx}{m + lx}. \end{aligned}$$

Hence, it follows in particular that the random variable ρ in this case does not have any moments of the first and higher orders, as distinct from the situation described in the previous section. However, some other characteristics of Bayesian queueing systems, depending on random variable $\rho = \lambda/\mu$, can have finite moments. Let us consider, for example, the queueing system of M|M|1|0 type. The probability that a claim has been received by the system will not be lost in a steady-state mode is equal to $\pi = 1/(1 + \rho)$ according to Erlangian formulae. As for the Bayesian problem statement, this probability becomes “random” by itself. Let us consider the distribution of the random variable π under the conditions of the model under study.

Provided $0 \leq y \leq 1$

$$P\{\pi < y\} = P\{\rho > (1 - y)/y\} = \frac{my}{my + l(1 - y)}$$

Correspondingly, the random variable π density is equal to $\frac{ml}{[my + l(1 - y)]^2}$, while the averaged probability that the call is not lost looks as follows

$$\mathbf{E}\pi = \int_0^1 \frac{mly}{[my + l(1 - y)]^2} dy = \frac{ml}{(m - l)^2} \left(\ln \frac{m}{l} + \frac{l}{m} - 1 \right).$$

It would be easy to calculate also the second moment of the random variable π as well as its other characteristics. Let us note that for $m = l$

$$\mathbf{E}\pi = 1/2.$$

The value

$$\pi = 1/(1 + \rho) = \frac{\mu}{\lambda + \mu}$$

is equal to value of the availability factor k (see above). Hence, the distribution of the random availability factor in case of exponentially distributed λ and μ is presented above as the distribution of random value π .

4. Erlang model for the parameter of service

In this chapter, the results will be presented for Erlang model of service time distribution.

Let us consider the system $M|M|1|0$ again. Let the parameter of input flow λ be degenerated, and the parameter of service μ has Erlang distribution with parameters n and α . First of all we will obtain the distribution functions and densities of random variables

$$\rho = \frac{\lambda}{\mu} \quad \text{and} \quad k = \pi = \frac{1}{1 + \rho}.$$

Let us find the distribution function $F_\rho(x)$ of the random variable ρ . We have

$$\begin{aligned} F_\rho(x) &= 1 - P\left(\mu < \frac{\lambda}{x}\right) = 1 - \int_0^{\lambda/x} \frac{t^{n-1} \alpha^n e^{-\alpha t}}{(n-1)!} dt \\ &= 1 - \frac{1}{(n-1)!} \int_0^{\alpha\lambda/x} z^{n-1} e^{-z} dz = e^{-\frac{\alpha\lambda}{x}} \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k}{x^k k!}, \quad x > 0. \end{aligned}$$

Having differentiated the last equation by x we obtain the density of variable ρ :

$$f_\rho(x) = e^{-\frac{\alpha\lambda}{x}} \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k (\alpha\lambda - kx)}{k! x^{k+2}}, \quad x > 0.$$

It is evident that within the described model the random variable ρ does not have moments of the first and the following rates:

$$\mathbf{E}\rho = \int_0^\infty e^{-\frac{\alpha\lambda}{x}} \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k (\alpha\lambda - kx)}{k! x^{k+1}} dx = \infty.$$

Consider probabilistic characteristics of the probability of “not losing” the claim, i. e. π . For the distribution function we have

$$F_{\pi}(x) = 1 - \mathbf{P}\left(\rho < \frac{1-x}{x}\right) = 1 - e^{-\frac{\alpha\lambda x}{1-x}} \cdot \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k x^k}{(1-x)^k k!},$$

$$x \in (0, 1).$$

In this case the density of π can be found with the formula

$$\begin{aligned} f_{\pi}(x) &= e^{-\frac{\alpha\lambda x}{1-x}} \frac{\alpha\lambda}{(1-x)^2} \cdot \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k x^k}{(1-x)^k k!} \\ &\quad - e^{-\frac{\alpha\lambda x}{1-x}} \cdot \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k}{k!} \cdot \frac{kx^{k-1}(1-x)^k + kx^k(1-x)^{k-1}}{(1-x)^{2k}} \\ &= e^{-\frac{\alpha\lambda x}{1-x}} \cdot \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k x^{k-1}(\alpha\lambda x - kx + k)}{k! (1-x)^{k+2}}, \quad x \in (0, 1). \end{aligned}$$

Let us find the expectation and the variance of the random variable π . We have

$$\begin{aligned} \mathbf{E}\pi &= \int_0^1 e^{-\frac{\alpha\lambda x}{1-x}} \cdot \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k x^k (\alpha\lambda x - kx + k)}{k! (1-x)^{k+2}} dx \\ &= \int_0^{\infty} e^{-z} \cdot \sum_{k=0}^{n-1} \frac{z^k (z + k)}{k! (\alpha\lambda + z)} dz \\ &= \sum_{k=0}^{n-1} \frac{1}{k!} \left[\int_0^{\infty} \frac{e^{-z} z^{k+1}}{\alpha\lambda + z} dz + k \int_0^{\infty} \frac{e^{-z} z^k}{\alpha\lambda + z} dz \right]. \end{aligned}$$

Consider the following designation. Let $\text{Ei}(x)$ be an integral exponential function:

$$\text{Ei}(x) = - \int_{-x}^{\infty} \frac{e^{-t}}{t} dt.$$

Using formulae [4, formula 567.9] and [5, formula 3.351], calculate as $k \geq 1$ the integral

$$\begin{aligned} \int_0^\infty \frac{e^{-z} z^k}{\alpha\lambda + z} dz &= \int_{\alpha\lambda}^\infty \frac{e^{\alpha\lambda} e^{-t} (t - \alpha\lambda)^k}{t} dt \\ &= e^{\alpha\lambda} (-\alpha\lambda)^k \int_{\alpha\lambda}^\infty \frac{e^{-t}}{t} dt + e^{\alpha\lambda} \int_{\alpha\lambda}^\infty e^{-t} \cdot \sum_{l=1}^k C_k^l t^{l-1} (-\alpha\lambda)^{k-l} dt \\ &= -e^{\alpha\lambda} (-\alpha\lambda)^k \text{Ei}(-\alpha\lambda) + \sum_{l=1}^k \sum_{m=0}^{l-1} (-1)^{k-l} (\alpha\lambda)^{k-l+m} \frac{k!}{l(k-l)! m!}. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbf{E}\pi &= 1 + e^{\alpha\lambda} \alpha\lambda \text{Ei}(-\alpha\lambda) - \sum_{k=1}^{n-1} \frac{e^{\alpha\lambda} (-\alpha\lambda)^k \text{Ei}(-\alpha\lambda) (k - \alpha\lambda)}{k!} \\ &\quad + \sum_{k=1}^{n-1} \sum_{l=1}^{k+1} \sum_{m=0}^{l-1} (-1)^{k-l+1} (\alpha\lambda)^{k-l+m+1} \frac{k+1}{l(k-l+1)! m!} \\ &\quad + \sum_{k=1}^{n-1} \sum_{l=1}^k \sum_{m=0}^{l-1} (-1)^{k-l} (\alpha\lambda)^{k-l+m} \frac{k}{l(k-l)! m!}. \end{aligned}$$

Now let us find the second moment of the random variable π :

$$\begin{aligned} \mathbf{E}\pi^2 &= \int_0^1 e^{-\frac{\alpha\lambda x}{1-x}} \sum_{k=0}^{n-1} \frac{(\alpha\lambda)^k x^{k+1} (\alpha\lambda x - kx + k)}{k! (1-x)^{k+2}} dx \\ &= \sum_{k=0}^{n-1} \frac{1}{k!} \left[\int_0^\infty \frac{e^{-z} z^{k+2}}{(\alpha\lambda + z)^2} dz + k \int_0^\infty \frac{e^{-z} z^{k+1}}{(\alpha\lambda + z)^2} dz \right]. \end{aligned}$$

Calculate as $k \geq 1$ the integral

$$\int_0^\infty \frac{e^{-z} z^{k+2}}{(\alpha\lambda + z)^2} dz = \int_{\alpha\lambda}^\infty \frac{e^{\alpha\lambda} e^{-t} (t - \alpha\lambda)^{k+2}}{t^2} dt$$

$$\begin{aligned}
&= e^{\alpha\lambda} \left[\int_{\alpha\lambda}^{\infty} \frac{e^{-t}(-\alpha\lambda)^{k+2}}{t^2} dt + \int_{\alpha\lambda}^{\infty} \frac{e^{-t}(k+2)(-\alpha\lambda)^{k+1}}{t} dt \right. \\
&\quad \left. + \sum_{l=2}^{k+2} C_{k+2}^l (-\alpha\lambda)^{k-l+2} \int_{\alpha\lambda}^{\infty} e^{-t} t^{l-2} dt \right] = (-1)^{k+2} (\alpha\lambda)^{k+1} \\
&\quad + e^{\alpha\lambda} (-\alpha\lambda)^{k+2} \text{Ei}(-\alpha\lambda) - e^{\alpha\lambda} (k+2)(-\alpha\lambda)^{k+1} \text{Ei}(-\alpha\lambda) \\
&\quad + \sum_{l=2}^{k+2} C_{k+2}^l (-\alpha\lambda)^{k-l+2} (l-2)! \sum_{m=0}^{l-2} \frac{(\alpha\lambda)^m}{m!}.
\end{aligned}$$

Thus,

$$\begin{aligned}
\mathbf{E}\pi^2 &= \sum_{k=0}^{n-1} \frac{1}{k!} \left[(-1)^{k+1} (\alpha\lambda)^k (k - \alpha\lambda) \right. \\
&\quad + e^{\alpha\lambda} (-\alpha\lambda)^k \text{Ei}(-\alpha\lambda) ((\alpha\lambda)^2 + 2\alpha\lambda - k(k+1)) \\
&\quad + \sum_{l=2}^{k+2} C_{k+2}^l (-\alpha\lambda)^{k-l+2} (l-2)! \sum_{m=0}^{l-2} \frac{(\alpha\lambda)^m}{m!} \\
&\quad \left. + k \sum_{l=2}^{k+1} C_{k+1}^l (-\alpha\lambda)^{k-l+1} (l-2)! \sum_{m=0}^{l-2} \frac{(\alpha\lambda)^m}{m!} \right].
\end{aligned}$$

It is now easy to find the variance of the random variable π .

Notice that despite cumbersome view of the formulae for the moments of the probability of “not losing” a claim, it is easy to realize them by computer to obtain the numerical expression of them for applied tasks for any natural n . However, the most interesting case of this model corresponds to $n = 1$ (i. e. exponential case). Let us quote the corresponding formulae:

$$F_{\rho}(x) = e^{-\frac{\alpha\lambda}{x}}, \quad f_{\rho}(x) = \frac{\alpha\lambda e^{-\frac{\alpha\lambda}{x}}}{x^2}, \quad x > 0;$$

$$\mathbf{E}\rho = \int_0^{\infty} \frac{\alpha\lambda e^{-\alpha\lambda y}}{y} dy = \infty;$$

$$F_{\pi}(x) = 1 - e^{-\frac{\alpha\lambda x}{1-x}}, \quad f_{\pi}(x) = \frac{\alpha\lambda e^{-\frac{\alpha\lambda x}{1-x}}}{(1-x)^2}, \quad x \in (0, 1);$$

$$\mathbf{E}\pi = 1 + \alpha\lambda e^{\alpha\lambda} \text{Ei}(-\alpha\lambda), \quad \mathbf{E}\pi^2 = \mathbf{E}\pi(2 + \alpha\lambda) - 1.$$

The results presented in this article are not complete within the problem of Bayesian queueing systems and Bayesian reliability problems; even considering $M|M|1|0$ type of systems. Obviously, further development within the presented problem area requires us to examine other a priori distributions of the variables λ , μ and other traditional input parameters of queueing models that could be interesting in practical cases. The distributions of the variables that characterize the functioning of different system types (such as $M|G|1$, $M|M|n|0$ and others) can be calculated after they have been randomized taking into account the given a priori distributions.

References

1. *Shorgin S. Ya.* On Bayesian queueing models // The 2nd Scientific Session of the Institute of Informatics Problems of the Russian Academy of Sciences: Reports theses. Moscow, IPI RAN, 2005. P. 120–121 (in Russian).
2. *Kozlov B. A., Ushakov I. A.* Reliability Handbook. — Holt, Rinehart & Winston, 1970.
3. *D'Apice C., Manzo R., Shorgin S.* Some Bayesian queueing and reliability models // Electronic J. Reliability: Theory & Applications. December 2006. V. 1, No. 4.
4. *Dwight H. B.* Tables of Integrals and Other Mathematical Data. — N. Y., The Macmillan Company, 1961.
5. *Gradshteyn I. S. and Ryzhik I. M.* Tablitsy integralov, summ, ryadov i proizvedenii (Tables of Integrals, Sums, Series, and Products). — Moscow, Nauka Academic Press, 1971 (in Russian).

G-NETWORK WITH THE ROUTE CHANGE ¹

Rosanna Manzo

Department of Information Engineering and Applied
Mathematics University of Salerno, Fisciano, Italy

Alexander Pechinkin

Institute of Informatics Problems, Russian Academy
of Sciences, Moscow, Russia

Queueing networks with negative customers (G-networks), Poisson flow of positive customers, non-exponential nodes, and dependent service at the different nodes are under consideration. Every customer arriving at the network is defined by a set of random parameters: customer route, the length of customer route, customer volume and its service time at each route stage as well. The arrival of a negative customer to a queueing system causes one of the ordinary (or “positive”) customers to be removed (or “killed”) if any is present. The “killed” customer continues its way along the new random route. For such G-networks, the multidimensional stationary distribution of the network state probabilities is shown to be representable in product form.

1. Network description

Recently, the big attention in the queueing theory is given to queueing networks with negative customers, or to the G-networks introduced by E. Gelenbe (see, for example, [1, 4]).

This attention is caused by the fact that G-networks model many phenomena in information and telecommunication networks, neural networks etc. You can find the extensive bibliography on G-networks in [5].

In the present paper, the variant of a G-network with dependent service at nodes is considered. Some previous results on G-networks with dependent service at nodes could be found in [6, 7]. Our new results, in terms of information and telecommunication networks, can be treated, in particular, as destruction of information message and their transformation in interfering messages (e. g. spam) which create additional loading for a network.

¹ The research has been performed with the support of the Russian Foundation for Basic Research (grants 06-07-89056 and 08-07-00152).

We consider an open queueing network with M nodes.

Each node s , $s = \overline{1, M}$, can be any of the following types:

- (1) infinite-server;
- (2) single-server with infinite buffer and LIFO discipline with interruption and resumption;
- (3) single-server with infinite buffer and processor sharing (PS) discipline.

Let's agree to denote random variables (RV) with capital Latin letters, and realizations of RV with corresponding lower case letters. Additionally, vector RV and any vectors we shall allocate with a semiboldface font.

A Poisson flow of (usual, positive) customers of intensity λ enters the network. Each customer arriving at the network is characterized by a set of random variables $(L, \mathbf{R}, \mathbf{Y}, \mathbf{X})$, which depend neither on analogous random variables for other customers nor on network pre-history, where:

- L is a customer route random length, i.e. the number of stages (nodes) at which he will be served;
- $\mathbf{R} = (R_1, \dots, R_L)$ is a random route comprising an assembly of node numbers (the same nodes at different stages are allowed) that the customer passes through in consecutive order at all L stages;
- $\mathbf{Y} = (Y_1, \dots, Y_L)$ are customer random volumes at route stages the customer consecutively passes through (the case when these volumes are different at different stages are considered too);
- $\mathbf{X} = (X_1, \dots, X_L)$ are customer random service times at the route stages the customer consecutively passes through.

It is obvious that under this network description the volume Y_n and the service time X_n define the service of a customer at node R_n . Let us recall that the routes \mathbf{R} with repetitions of node numbers are allowed, i.e. a customer can be served at the same node s several times (but probably with different customer volumes).

Stochastic characteristics of a random variable $(L, \mathbf{R}, \mathbf{Y}, \mathbf{X})$ are given by the joint probability distribution function (PDF)

$$B(l, \mathbf{r}, \mathbf{y}, \mathbf{x}) = \mathbf{P}\{L = l, R_n = r_n, Y_n \leq y_n, X_n \leq x_n, \quad n = \overline{1, l}\}.$$

Further on, let us denote by

$$G(l, \mathbf{r}, \mathbf{y}) = \mathbf{P}\{L = l, R_n = r_n, Y_n \leq y_n, \quad n = \overline{1, l}\}$$

the joint PDF of the route \mathbf{R} of the length L and with customer volumes \mathbf{Y} at the route stages, by

$$B(\mathbf{x} | l, \mathbf{r}, \mathbf{y}) = \mathbf{P}\{X_n \leq x_n, \quad n = \overline{1, l} \mid L = l, \mathbf{R} = \mathbf{r}, \mathbf{Y} = \mathbf{y}\}$$

the conditional joint PDF of the customer service lengths \mathbf{X} at the route stages under a fixed route $\mathbf{R} = \mathbf{r}$ of the length $L = l$ and volumes $\mathbf{Y} = \mathbf{y}$, and by

$$B_n(x | l, \mathbf{r}, \mathbf{y}) = \mathbf{P}\{X_n \leq x \mid L = l, \mathbf{R} = \mathbf{r}, \mathbf{Y} = \mathbf{y}\}, \quad n = \overline{1, l},$$

the conditional PDF of the customer service length X_n at the n -th stage (at a node with number $R_n = r_n$) under a fixed route $\mathbf{R} = \mathbf{r}$ of the length $L = l$ and volumes $\mathbf{Y} = \mathbf{y}$.

We shall expect that service lengths are conditionally independent along the route, i. e. the conditional PDF $B(\mathbf{x} | l, \mathbf{r}, \mathbf{y})$ has the form

$$B(\mathbf{x} | l, \mathbf{r}, \mathbf{y}) = \prod_{n=1}^l B_n(x_n | l, \mathbf{r}, \mathbf{y}).$$

Along with the flow of positive customers described above, flows of negative customers arrive at the network. These flows are defined in the following way.

A1. The flows arriving at different nodes are independent.

A2. A customer flow arriving at node s of type 2 or 3 is Poisson one of intensity γ_s .

A3. A customer flow arriving at node s of type 1 is a Markov one with intensity $\gamma_s(n)$ depending (only) on the number n of busy servers at this node in the following way: $\gamma_s(n) = ns$.

A4. A negative customer arriving at a node s with k positive customers in service at this node (if this node is of the type 1 or 3, then k is the total number of positive customers at the node, if it is of type 2, then $k = 1$) chooses one of positive customers being served with probability $1/k$. After this, if the chosen customer is “not killed”, the negative customer either with probability $\omega_n(x | l, \mathbf{r}, \mathbf{y})$ immediately “kills” it and quits the network or with the complementary probability $1 - \omega_n(x | l, \mathbf{r}, \mathbf{y})$ quits the network without inducing any action. Here $(l, \mathbf{r}, \mathbf{y})$ are the parameters of the chosen positive customer, defined earlier; n is the number of route stage in which this customer is served (it is only natural that $r_n = s$) and x is the elaborated (served) customer length. However, the “killed” positive customer does not leave the system

but passes in node R_1^* and continues to be served according to new RV $(L^*, \mathbf{R}^*, \mathbf{Y}^*, \mathbf{X}^*)$. Distribution of RV $(L^*, \mathbf{R}^*, \mathbf{Y}^*, \mathbf{X}^*)$ depends only on RV $(L, \mathbf{R}, \mathbf{Y})$ of “killed” customer and N , i.e., number of stage at which a customer has been “killed” in node $R_N = s$. The distribution has, under condition of $(L, \mathbf{R}, \mathbf{Y}) = (l, \mathbf{r}, \mathbf{y})$ and when the customer has been “killed” at the stage with number $N = n$, a conditional distribution function

$$H(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n) \\ = \mathbf{P}\{L^* = l^*, R_m^* = r_m^*, Y_m^* \leq y_m^*, m = \overline{1, l^*} | \\ (L, \mathbf{R}, \mathbf{Y}) = (l, \mathbf{r}, \mathbf{y}), N = n\}.$$

Service lengths of “killed” customer are depends from $(L^*, \mathbf{R}^*, \mathbf{Y}^*)$ only and are conditionally independent along the route, i.e. the conditional PDF

$$C(\mathbf{x} | l^*, \mathbf{r}^*, \mathbf{y}^*) = \mathbf{P}\{X_m^* \leq x_m, m = \overline{1, l^*} | \\ L^* = l^*, \mathbf{R}^* = \mathbf{r}^*, \mathbf{Y}^* = \mathbf{y}^*\}$$

is of the form of

$$C(\mathbf{x} | l^*, \mathbf{r}^*, \mathbf{y}^*) = \prod_{m=1}^{l^*} C_m(x_m | l^*, \mathbf{r}^*, \mathbf{y}^*),$$

where

$$C_m(x | l^*, \mathbf{r}^*, \mathbf{y}^*) = \mathbf{P}\{X_m^* \leq x | L^* = l^*, \mathbf{R}^* = \mathbf{r}^*, \mathbf{Y}^* = \mathbf{y}^*\}, \\ m = \overline{1, l^*},$$

is the conditional PDF of the “killed” customer service length X_n^* at the m -th stage (at a node with number $R_n^* = r_n^*$) under a fixed route $\mathbf{R}^* = \mathbf{r}^*$ of the length $L^* = l^*$ and volumes $\mathbf{Y}^* = \mathbf{y}^*$. If there is a “killed” already positive customer on the chosen server, it simply passes to the node corresponding to the consequent stage of service and continues to be served accordingly to its set of parameters. Finally, if at the moment of a negative customer’s arrival into some node there is no positive customers there, then the negative customer quits the network without inducing any action.

We shall make an additional technical assumption on PDFs $G(l, \mathbf{r}, \mathbf{y})$ and $H(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n)$. Namely, we suppose that the

PDFs $G(l, \mathbf{r}, \mathbf{y})$ and $H(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n)$ are absolutely continuous, and denote by $g(l, \mathbf{r}, \mathbf{y})$ and $h(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n)$ their densities, i. e.

$$g(l, \mathbf{r}, \mathbf{y}) = \frac{\partial^l}{\partial y_1 \dots \partial y_l} G(l, \mathbf{r}, \mathbf{y}).$$

and

$$h(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n) = \frac{\partial^{l^*}}{\partial y_1^* \dots \partial y_{l^*}^*} H(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n).$$

This assumption could be easily neglected if we interpret derivatives as generalized ones.

2. Auxiliary functions

First we introduce auxiliary functions in the following way (below we shall use the notation $\bar{\omega} = 1 - \omega$ for any probability ω , in particular, $\bar{F}(x) = 1 - F(x)$ for any PDF $F(x)$):

$$\begin{aligned} F_n(x | l, \mathbf{r}, \mathbf{y}) &= 1 - \exp \left\{ -\gamma_{r_n} \int_0^x \omega_n(z | l, \mathbf{r}, \mathbf{y}) dz \right\}, \quad n = \overline{1, l}, \\ B_n^*(x | l, \mathbf{r}, \mathbf{y}) &= 1 - \bar{B}_n(x | l, \mathbf{r}, \mathbf{y}) \bar{F}_n(x | l, \mathbf{r}, \mathbf{y}), \quad n = \overline{1, l}, \\ \omega_n(l, \mathbf{r}, \mathbf{y}) &= \int_0^\infty \bar{F}_n(x | l, \mathbf{r}, \mathbf{y}) b_n(x | l, \mathbf{r}, \mathbf{y}) dx, \quad n = \overline{1, l}, \\ \omega_n^*(l, \mathbf{r}, \mathbf{y}) &= \prod_{i=1}^{n-1} \omega_i(l, \mathbf{r}, \mathbf{y}), \quad n = \overline{1, l+1}, \\ g_n^*(l, \mathbf{r}, \mathbf{y}) &= \omega_n^*(l, \mathbf{r}, \mathbf{y}) g(l, \mathbf{r}, \mathbf{y}), \quad n = \overline{1, l}, \\ h(l^*, \mathbf{r}^*, \mathbf{y}^*) &= \sum_{l, \mathbf{r}} \int_{\mathbb{R}^l} \sum_{n=1}^l h(l^*, \mathbf{r}^*, \mathbf{y}^* | l, \mathbf{r}, \mathbf{y}, n) g_n^*(l, \mathbf{r}, \mathbf{y}) \bar{\omega}_n(l, \mathbf{r}, \mathbf{y}) d\mathbf{y}, \\ m_n^+(l, \mathbf{r}, \mathbf{y}) &= \int_0^\infty \bar{B}_n^*(x | l, \mathbf{r}, \mathbf{y}) dx, \quad n = \overline{1, l}, \\ m_n^-(l, \mathbf{r}, \mathbf{y}) &= \int_0^\infty \bar{C}_n(x | l, \mathbf{r}, \mathbf{y}) dx, \quad n = \overline{1, l}. \end{aligned}$$

Hereinafter summation by vector argument means summation by all possible values of its coordinates. Besides, for the sake of brevity we shall use the notations

$$\int_{\mathbb{R}^l} \dots d\mathbf{y} = \int_{\mathbb{R}^l} \dots dy_1 \dots dy_l.$$

It is only natural that for the nodes of types 2 and 3 the last two characteristics are defined under the condition that there exist no other customers at these nodes.

Let us set for the s -th node

$$\lambda_s^+ = \lambda \sum_{l, \mathbf{r}} \int \sum_{\mathbb{R}^l}^l \delta_{s-r_n} g_n^*(l, \mathbf{r}, \mathbf{y}) d\mathbf{y},$$

$$\lambda_s^- = \lambda \sum_{l, \mathbf{r}} \int \sum_{\mathbb{R}^l}^l \delta_{s-r_n} h(l, \mathbf{r}, \mathbf{y}) d\mathbf{y},$$

$$\lambda_s = \lambda_s^+ + \lambda_s^-,$$

$$\rho_s^+ = \lambda \sum_{l, \mathbf{r}} \int \sum_{\mathbb{R}^l}^l \delta_{s-r_n} g_n^*(l, \mathbf{r}, \mathbf{y}) m_n^+(l, \mathbf{r}, \mathbf{y}) d\mathbf{y},$$

$$\rho_s^- = \lambda \sum_{l, \mathbf{r}} \int \sum_{\mathbb{R}^l}^l \delta_{s-r_n} h(l, \mathbf{r}, \mathbf{y}) m_n^-(l, \mathbf{r}, \mathbf{y}) d\mathbf{y},$$

$$\rho_s = \rho_s^+ + \rho_s^-,$$

where δ_j is the Kronecker symbol.

Let us suppose that $\lambda_s < \infty$ for all nodes s . The last condition means that the total flow intensity λ_s of positive customers arriving at node s is finite. Note that this condition does not follow from the condition that the traffic intensity of these nodes is finite too (the latter condition will be given below).

3. Markov process

Let us now define the Markov process describing the stochastic behavior of the network under consideration.

We shall denote a network state by an assembly $\mathbf{z} = (z_1, \dots, z_M)$, where the assembly $\mathbf{z}_s = (k_s, z_{s1}, \dots, z_{sk_s})$, $s = \overline{1, M}$, in turn, de-

scribes the state of the s th node in the following way: k_s is the number of customer at the s th node and the assembly \mathbf{z}_{si} , $s = \overline{1, M}$, $i = \overline{1, k_s}$, with components $\mathbf{z}_{si} = (l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}, w_{si}, n_{si}, x_{si})$ stores the information $(l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}, w_{si})$ on the i -th customer at the s th node, and its position (n_{si}, x_{si}) in the network:

- l_{si} is the route length;
- $\mathbf{r}_{si} = (r_{si1}, \dots, r_{sil_{si}})$ is the route;
- $\mathbf{y}_{si} = (y_{si1}, \dots, y_{sil_{si}})$ are customer volumes at its route stages;
- n_{si} is the number of the route's stage at which the customer exists (while being served or waiting for service); clearly, $n_{si} \leq l_{si}$;
- w_{si} is the function which shows customer state; we set $w_{si} = 0$ if the customer is “not killed”, and $w_{si} = 1$ if the customer is “killed” (but is being served);
- x_{si} is the customer length already serviced at a given stage.

Evidently that due to the notations introduced above, we have $r_{sin_{si}} = s$. It is also clear that the vector $\mathbf{z}_s = \mathbf{0}$ if $k_s = 0$, i.e. when there are no customers at the s th node, and the vector $\mathbf{z} = \mathbf{0} = (0, \dots, 0)$ in the case, when there are no customers in the network at all.

In what follows we will accept the following rule of numbering of customers in the nodes. For the nodes of types 1 or 3, the numbers are assigned randomly, and for the nodes of type 2, in the inverse order to their arrivals at the nodes.

The set of states of the network is denoted by $\mathcal{Z} = \{\mathbf{z}\}$.

To describe the operation of our queueing network, let us consider the process

$$\mathbf{Z}(t) = \mathbf{z}, \text{ if the network exists in state } \mathbf{z} \text{ at instant } t.$$

It is obviously a Markov process.

4. Product form solution

The stationary density of the state probability distribution of the process $\mathbf{Z}(t)$ is denoted by $p(\mathbf{z})$.

The following theorem on the multiplicative representation of stationary network state probabilities for a considered network holds.

Theorem. If for a node s of type 1 $\rho_s < 1$, and for a node s of types 2 and 3 $\rho_s < 1$, then there exists a limit (stationary) probability state distribution of the process $\mathbf{Z}(t)$ with probability distribution density

$$p(\mathbf{z}) = \prod_{s=1}^M p_s(\mathbf{z}_s),$$

thereby:

for a node s of the type 1

$$p_s(\mathbf{z}_s) = e^{-\rho_s} \frac{\lambda_s^{k_s}}{k_s!} \prod_{i=1}^{k_s} (\delta_{w_{si}} g_{n_{si}}^*(l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) \bar{B}_{n_{si}}^*(x_{si} | l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) + \delta_{1-w_{si}} h(l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) \bar{C}_{n_{si}}(x_{si} | l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}));$$

for a node s of the type 2 or 3

$$p_s(\mathbf{z}_s) = (1 - \rho_s) \lambda_s^{k_s} \prod_{i=1}^{k_s} (\delta_{w_{si}} g_{n_{si}}^*(l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) \bar{B}_{n_{si}}^*(x_{si} | l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) + \delta_{1-w_{si}} h(l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si}) \bar{C}_{n_{si}}(x_{si} | l_{si}, \mathbf{r}_{si}, \mathbf{y}_{si})).$$

References

1. Gelenbe E. Random neural networks with positive and negative signals and product form solution // Neural Computation. 1989. V.:1. P.:502–510.
2. Gelenbe E. Reseaux stochastiques ouverts avec clients negatifs et positifs, et reseaux neuronaux // Comptes-Rendus Acad. Sci. Paris II. 1989. V.:309. P.:979–982.
3. Gelenbe E. Reseaux neuronaux aleatoires stables // Comptes-Rendus Acad. Sci. II. 1990. V.:310. P.:177–180.
4. Gelenbe E. Stable random neural networks // Neural Computation. 1990. V.:2. P.:239–247.
5. Bocharov P., P. and Vishnevskii V., M. G-networks: development of the theory of multiplicative networks // Automation and Remote Control. 2003. V.:64. P.:714–739.
6. Bocharov P., P., D'Apice C., Gavrilov E., V. and Pechinkin A., V. Product form solution for G-networks with dependent service // RAIRO Oper. Res. 2004. V.:38. P.:105–119.
7. Bocharov P., P., D'Apice C. and Pechinkin A., V. Product form solution for exponential G-networks with dependent service and completion of service of “killed” customers // Computational Management Science. 2006. No. 3. P. 177–192.

FAN-BEAM STOCHASTIC TOMOGRAPHY¹

Oleg Shestakov

Faculty of Computational Mathematics and Cybernetics,
Lomonosov Moscow State University, Moscow, Russia

In fan-beam tomography settings an object is illuminated by the divergent pencil of rays emitted from the source which moves around it. This scheme works much faster than traditional parallel beam tomography. In some biological and physical applications an object under study is described by random function. This paper considers the problem of recovering probabilistic characteristics of random function from fan-beam projections.

1. Introduction. Problems of recovering probabilistic characteristics of bivariate random functions from characteristics of univariate projections arise in physics and microbiology (see [1]). The main feature of these problems is that object under study may have several (and even infinite number) of states, which randomly change during the process of projection registration and conventional tomographic approach cannot be used in this case. Such problems led to appearance of a new branch of computational tomography, which was named stochastic tomography.

In papers [2-4] the problem of recovering probabilistic characteristics of bivariate random functions is considered in parallel beam settings. It is shown that in general case this problem is characterized by a great ambiguity and without restrictions on realizations of random functions meaningful results can be obtained only in the case, when random function has at most denumerable number of states. Distribution reconstruction method for class of such functions is introduced in paper [4].

Sometimes, however, one must use fan-beam scheme (see [1]), where an object is illuminated by the divergent pencil of rays emitted from the source which moves around it. In this paper we consider the problem of recovering probabilistic characteristics of bivariate random functions from characteristics of univariate projections for such scheme.

¹ This work was supported by the Russian Foundation for Basic Research (grant 08-01-00567-a).

2. Statement of the problem. A formal statement of the problem is as follows. There is a bivariate random function $\xi(x, y)$, which describes stochastic object. Throughout this paper we will assume that: 1) $\xi(x, y)$ has a compact support (without loss of generality we will assume that this support is a unit circle with center at the origin: $U = \{(x, y) \in \mathbf{R}^2 : x^2 + y^2 \leq 1\}$), 2) $\xi(x, y)$ is integrable with probability 1: $\xi(x, y) \in L^1(U)$ (a. e.). Functions, which differ on a set of null Lebesgue measure, will be considered as equivalent. Thus the sign “ \equiv ” means equivalence in L^1 -norm.

Let us assume that the source of radiation moves around the circle with radius D centered at the origin and the circle U is contained within this circle. Denote by $R_\beta \xi(\gamma)$ projection of function $\xi(x, y)$ in fan-beam scheme of scanning. Here $\beta \in [0, 2\pi)$ is an angular coordinate of the source (an angle between vertical axis and a line connecting origin and the source of radiation), and $\gamma \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$ describes the position of ray within the pencil of rays (an angle between the given ray and the central ray of the pencil). In this scheme projection data are described by the following expression:

$$R_\beta \xi(\gamma) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \xi(x, y) \times \delta(x \cos(\beta + \gamma) + y \sin(\beta + \gamma) - D \sin \gamma) dx dy.$$

In stochastic tomography problems it is supposed that some information on probabilistic characteristics of projections (all or some set) is available. The problem is to determine certain probabilistic characteristics of the random function $\xi(x, y)$.

The first question is whether this determination is possible, i. e. about uniqueness of the correspondence between characteristics of bivariate random function and characteristics of its projections. As in parallel beam settings this question may take the following forms:

1. Is it possible to unambiguously determine joint distributions of $\xi(x_1, y_1), \dots, \xi(x_n, y_n)$, if joint distributions of $R_\beta \xi(\gamma_1), \dots, R_\beta \xi(\gamma_m)$ are known for all $m = 1, 2, \dots$ and all $\beta \in [0, 2\pi)$?

2. Is it possible to unambiguously determine variances $D(\xi(x, y))$, if variances $D(R_\beta \xi(\gamma))$ are known for all $\beta \in [0, 2\pi]$ and $\gamma \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$?

3. Is it possible to unambiguously determine variances $D(\xi(x, y))$, if joint distributions of $R_\beta \xi(\gamma_1), \dots, R_\beta \xi(\gamma_m)$ are known for all $m = 1, 2, \dots$ and all $\beta \in [0, 2\pi)$?

4. Is there any connection between “variability” (the variance value) of a random bivariate function in a given point and variability of projections in projections of this point (this question is important for biological applications)?

We do not consider the problem of recovering mathematical expectation from mathematical expectations of projections because this problem is equivalent to conventional (nonstochastic) tomography (see [2]).

3. Radon transform. In this section we will provide necessary information on Radon transform, which lies in the basis of parallel beam tomography algorithms. Radon transform of integrable function $f(x, y)$ is defined by

$$P_\varphi f(t) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \delta(x \cos \varphi + y \sin \varphi - t) dx dy$$

for $\varphi \in [0, \pi)$ and $t \in \mathbf{R}$.

Central slice theorem holds (see [5]):

$$\widehat{P_\varphi f}(\omega) = \widehat{f}(\omega \cos \varphi, \omega \sin \varphi),$$

where $\widehat{P_\varphi f}(\omega)$ is an univariate Fourier transform of $P_\varphi f(t)$ with respect to variable t , and $\widehat{f}(\omega_1, \omega_2)$ is a bivariate Fourier transform of $f(x, y)$. General projection theorem also holds (see [5]):

$$\int_{-\infty}^{\infty} P_\varphi f(t) h(t) dt = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) h(\omega \cos \varphi + \omega \sin \varphi) dx dy, \quad (1)$$

where $h(t)$ is an arbitrary function, such that integrals in the left-hand and right-hand sides of (1) exist.

General projection theorem gives the following useful corollary (see [5]):

$$\int_{-\infty}^{\infty} P_{\varphi} f(t) t^m dt = p_m(\varphi),$$

where $p_m(\varphi)$ is a polynomial in $\cos \varphi$ and $\sin \varphi$ of degree m .

The following relation between $P_{\varphi} f(t)$ and $R_{\beta} f(\gamma)$ holds:

$$R_{\beta} f(\gamma) = P_{\beta+\gamma} f(D \sin \gamma).$$

4. General case. Counterexamples. In this section we will prove some propositions, which show that in general cases answers to the questions 2 and 4 from section 2 are negative. It seems that answers to the questions 1 and 3 are also negative, although there are no counterexamples for fan-beam scheme yet (counterexamples for parallel beam scheme can be found in [2] and [3]). However, in the next section we will describe a class of random functions (sufficient for biological applications), for which it is possible to reconstruct probabilistic characteristics of random function, if probabilistic characteristics of some set of projections are known.

Proposition 1. *There exist two random functions $\xi_1(x, y)$ and $\xi_2(x, y)$, defined on the unit circle U , such that $D(\xi_1(x, y)) \neq D(\xi_2(x, y))$ for all x and y , such that $x^2 + y^2 < 1$, while $D(R_{\beta} \xi_1(\gamma)) = D(R_{\beta} \xi_2(\gamma))$ for all $\beta \in [0, 2\pi)$ and $\gamma \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$.*

Proof. Let $\eta_1(x, y)$ and $\eta_2(x, y)$ be two homogeneous and isotropic random fields such that $E\eta_1(x, y) = E\eta_2(x, y) = 0$, $D(\eta_1(x, y)) = D(\eta_2(x, y)) = 1$. Suppose that the covariance function $\rho_1(t)$ of $\eta_1(x, y)$ is positive and decreases for $t > 0$, and the covariance function $\rho_2(t)$ of $\eta_2(x, y)$ satisfies the equality $\rho_2(t) = \rho_1(ct)$, where $c > 1$. Let $\eta_0(x, y)$ be a random field of the form $\eta_0(x, y) = \nu f(x, y)$, where ν is a random variable independent of $\eta_2(x, y)$ and taking two values, -1 and 1, with probability 1/2 each, and $f(x, y)$ is a spherically symmetric function, positive for $x^2 + y^2 < 1$ and equal to zero for $x^2 + y^2 \geq 1$. Below we prove that this function can be chosen in such a way that all conditions of the proposition are satisfied.

Let us define

$$\xi_1(x, y) = \begin{cases} \eta_1(x, y) & \text{for } x^2 + y^2 \leq 1, \\ 0 & \text{otherwise} \end{cases}$$

and

$$\xi_2(x, y) = \begin{cases} \eta_1(x, y) + \eta_0(x, y) & \text{for } x^2 + y^2 \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Then for $x^2 + y^2 < 1$,

$$\begin{aligned} D(\xi_2(x, y)) &= D(\eta_1(x, y)) + D(\eta_0(x, y)) = 1 + f^2(x, y) \neq 1 \\ &= D(\xi_1(x, y)). \end{aligned}$$

Let us prove that there exists a spherically symmetric, positive (for $x^2 + y^2 < 1$) function $f(x, y)$ such that variances of all projections of $\xi_1(x, y)$ and $\xi_2(x, y)$ coincide.

Due to the symmetry, it is sufficient to consider the case $\beta = 0$. We have

$$\begin{aligned} D(R_0\xi_1(\gamma)) &= D(P_\gamma\xi_1(D\sin\gamma)) = D(P_0\xi_1(D\sin\gamma)) \\ &= D\left[\int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \xi_1(D\sin\gamma, y) dy\right] \\ &= E\left[\int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \xi_1(D\sin\gamma, y) dy\right]^2 \\ &= E\left[\int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \xi_1(D\sin\gamma, u) du \cdot \int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \xi_1(D\sin\gamma, v) dv\right] \\ &= E\left[\int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \int_{-\sqrt{1-(D\sin\gamma)^2}}^{\sqrt{1-(D\sin\gamma)^2}} \xi_1(D\sin\gamma, u) \xi_2(D\sin\gamma, v) du dv\right] \end{aligned}$$

$$\begin{aligned}
&= \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} E \xi_1(D \sin \gamma, u) \xi_2(D \sin \gamma, v) du dv \\
&= \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \rho_1(u - v) du dv
\end{aligned}$$

and, analogously,

$$\begin{aligned}
&D(R_0 \xi_2(\gamma)) \\
&= \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \rho_1(c(u - v)) du dv + [R_0 f(\gamma)]^2.
\end{aligned}$$

From lemma 2 of [2] it follows that there exists positive (for $x^2 + y^2 < 1$) spherically symmetric function $f(x, y)$ such that

$$P_\varphi f(t) = \left[\int_{-\sqrt{1-t^2}}^{\sqrt{1-t^2}} \int_{-\sqrt{1-t^2}}^{\sqrt{1-t^2}} (\rho_1(u - v) - \rho_1(c(u - v))) du dv \right]^{1/2}.$$

For this function we have

$$\begin{aligned}
R_0 f(\gamma) &= P_\gamma f(D \sin \gamma) = P_0 f(D \sin \gamma) \\
&= \left[\int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} \int_{-\sqrt{1-(D \sin \gamma)^2}}^{\sqrt{1-(D \sin \gamma)^2}} (\rho_1(u - v) - \rho_1(c(u - v))) du dv \right]^{1/2}.
\end{aligned}$$

Thus $D(R_0 \xi_1(\gamma)) = D(R_0 \xi_2(\gamma))$. This completes the proof.

For some biological problems it is sufficient to determine the points, in which bivariate random function has relatively large variance. In other words there is no need to obtain variance values. Thus there is a question of relationship between variance magnitude of certain points of random object and variance magnitude of projections of these points. The following propositions show that in general case there is no relationship between these characteristics.

Proposition 2. *There exists a random function $\xi(x, y)$, $x^2 + y^2 \leq 1$, such that all projections of its most variable point (having maximal variance) are the least variable points of the corresponding projections.*

Proof. Let η be a random variable such that

$$E\eta = 0, \quad D(\eta) = 1.$$

Let us define

$$\psi(x, y) = \begin{cases} 1 - 2\sqrt{x^2 + y^2} & \text{if } \sqrt{x^2 + y^2} \leq \frac{1}{2}, \\ -\frac{1}{2} & \text{if } \frac{1}{2} < \sqrt{x^2 + y^2} \leq 1, \\ 0 & \text{if } \sqrt{x^2 + y^2} > 1. \end{cases}$$

Consider the following random field

$$\xi(x, y) = \eta\psi(x, y).$$

We have

$$E\xi(x, y) = 0,$$

and

$$D(\xi(x, y)) = E\xi^2(x, y) = \psi^2(x, y),$$

so, the point $(0, 0)$ is the most variable point of the random field $\xi(x, y)$ ($\psi(0, 0) > |\psi(x, y)|$, hence $D(\xi(0, 0)) > D(\xi(x, y))$ for all $(x, y) \neq (0, 0)$). Let us prove that all projections of this point have null variance:

$$D(R_\beta\xi(0)) = 0.$$

Due to the symmetry, it suffices to prove this equality for one value of β , say, $\beta = 0$. We have

$$D(R_0\xi(0)) = E\left[\int_{-1}^1 \eta\psi(0, y) dy\right]^2 = E\eta^2 \cdot \left[\int_{-1}^1 \psi(0, y) dy\right]^2 = 0.$$

At the same time, $D(R_0\xi(\gamma)) \neq 0$ for all $\gamma \neq 0$ ($\gamma \in [-\frac{\pi}{2}, \frac{\pi}{2}]$), since $R_0\xi(\gamma) = P_\gamma\xi(D\sin\gamma)$ and, as it is shown in [2], $D(P_\varphi\xi(t)) \neq 0$, if $t \neq 0$.

This completes the proof of the proposition.

Proposition 3. *There exists a random function $\xi(x, y)$, $x^2 + y^2 \leq 1$, such that all projections of its least variable point*

(having minimal variance) are the most variable points of the corresponding projections.

The proof is analogous to that of the previous proposition.

5. Class T and the uniqueness theorem. As it follows from the previous section, in general case, i. e. when stochastic object is described by an arbitrary random function, no meaningful results can be obtained. However, situation can change under some restrictions on random functions considered. In this section we introduce a class of random functions which, on the one hand, is wide enough to describe many situations in applications, for instance, in microbiology, and on the other hand, guarantees uniqueness of correspondence between probabilistic distributions of stochastic objects and distributions of projections.

Let T be the set of all random functions $\xi(x, y)$ of the form

$$\xi(x, y) = f_\nu(x, y),$$

where $f_1(x, y), f_2(x, y), \dots$ — is a sequence of integrable functions defined on the unite circle $U = \{(x, y) \in \mathbf{R}^2: x^2 + y^2 \leq 1\}$, and ν is a random variable taking positive integer values.

Random functions from class T are nothing else than discrete random elements in the space $L^1(U)$, therefore their probabilistic structure is completely determined by the distribution, i. e. by the collection

$$(f_1(x, y), f_2(x, y), \dots; p_1, p_2, \dots),$$

where $p_i = Pr(\xi(x, y) = f_i(x, y)), i = 1, 2, \dots, \sum_{i=1}^{\infty} p_i = 1$. We denote the distribution of $\xi(x, y)$ by Pr_ξ .

It turns out that in the frames of introduced model distribution of a stochastic object is uniquely determined by distributions of its projections.

Theorem 1. *Let $\xi(x, y) \in T, \eta(x, y) \in T$, and*

$$Pr_{R_\beta \xi} = Pr_{R_\beta \eta}$$

for all $\beta \in \Lambda \subseteq [0, 2\pi)$, where Λ is an infinite set, then

$$Pr_\xi = Pr_\eta.$$

In other words, in class T distribution of any random function is uniquely determined by distributions of each infinite set of its projections.

To prove this theorem we will make use of the following lemma, stating uniqueness of reconstruction from fan-beam projections (see [6]).

Lemma 1. *Let $f(x, y) \in L^1(U)$ and let $\Lambda \subseteq [0, 2\pi)$ be an infinite set. If $R_\beta f(\gamma) \equiv 0$ for all $\beta \in \Lambda$, then $f(x, y) \equiv 0$.*

Proof. Let us note that

$$R_\beta f(\gamma) = \int_{-\infty}^{\infty} f(-D \sin \beta + t \sin(\beta + \gamma), D \cos \beta - t \cos(\beta + \gamma)) dt.$$

Consider the following function

$$H(\beta, \varphi) = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \frac{R_\beta f(\gamma)}{|\sin(\beta + \gamma) \cos \varphi - \cos(\beta + \gamma) \sin \varphi|} d\gamma. \quad (2)$$

We have

$$\begin{aligned} H(\beta, \varphi) &= \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \int_{-\infty}^{\infty} \frac{f(-D \sin \beta + t \sin(\beta + \gamma), D \cos \beta - t \cos(\beta + \gamma)) |t|}{|t(\sin(\beta + \gamma) \cos \varphi - \cos(\beta + \gamma) \sin \varphi)|} dt d\gamma \\ &= \int_{-\frac{\pi}{2}+\beta}^{\frac{\pi}{2}+\beta} \int_{-\infty}^{\infty} \frac{f(-D \sin \beta + t \sin \vartheta, D \cos \beta - t \cos \vartheta) |t|}{|t(\sin \vartheta \cos \varphi - \cos \vartheta \sin \varphi)|} dt d\vartheta \\ &= \int_{-\frac{\pi}{2}+\beta}^{\frac{3\pi}{2}+\beta} \int_0^{\infty} \frac{f(-D \sin \beta + t \sin \vartheta, D \cos \beta - t \cos \vartheta) t}{|t(\sin \vartheta \cos \varphi - \cos \vartheta \sin \varphi)|} dt d\vartheta. \end{aligned}$$

Integrand is a 2π -periodic function of ϑ , so we can write

$$H(\beta, \varphi) = \int_0^{2\pi} \int_0^{\infty} \frac{f(-D \sin \beta + t \sin \vartheta, D \cos \beta - t \cos \vartheta) t}{|t(\sin \vartheta \cos \varphi - \cos \vartheta \sin \varphi)|} dt d\vartheta.$$

Switching from the polar to rectangular coordinate system and taking into consideration that $f(x, y)$ is defined in the unit circle U , we have

$$\begin{aligned}
 H(\beta, \varphi) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{f(x, y)}{|-D \sin \beta \cos \varphi + D \cos \beta \sin \varphi - x \cos \varphi - y \sin \varphi|} dx dy \\
 &= \int \int_U \frac{f(x, y)}{|D \sin(\varphi - \beta) - (x \cos \varphi + y \sin \varphi)|} dx dy \\
 &= \int_{-1}^1 \frac{P_\varphi f(t)}{|D \sin(\varphi - \beta) - t|} dt, \quad (3)
 \end{aligned}$$

where in the last equality we made use of general projection theorem (1). By the assumption of lemma $H(\beta, \varphi) = 0$ for $\beta \in \Lambda$. Let β_0 be any limit point of the set Λ , and let O be some fixed neighborhood of point $\beta_0 + \frac{\pi}{2}$, such that for any $\varphi \in O$ we have $D \sin(\varphi - \beta_0) > 1 + 2\varepsilon$. Then for an infinite number of $\beta \in \Lambda$ and any $\varphi \in O$ we have $D \sin(\varphi - \beta) > 1 + \varepsilon$. For these β and φ function $\frac{1}{D \sin(\varphi - \beta) - t}$ can be expanded into series, which converges uniformly for $t \in [-1, 1]$. Substituting this expansion into (3) and integrating termwise, we have

$$H(\beta, \varphi) = \sum_{m=0}^{\infty} (D \sin(\varphi - \beta))^{-(m+1)} p_m(\varphi). \quad (4)$$

Function $H(\beta, \varphi)$ vanishes on an infinite set of $\beta \in \Lambda$ and $\varphi \in O$, such that $D \sin(\varphi - \beta) > 1 + \varepsilon$. It is possible only if all $p_m(\varphi) = 0$ for any $\varphi \in O$. Since $p_m(\varphi)$ are polynomials, they are identically zero for all $\varphi \in [0, 2\pi]$. This means that $P_\varphi f(t) \equiv 0$ for all $\varphi \in [0, \pi)$, and hence by the central slice theorem we conclude that $f(x, y) \equiv 0$. This completes the proof of lemma.

Proof of the theorem 1. Let $\Lambda \subseteq [0, 2\pi)$ be an infinite set and $Pr_{R_\beta \xi} = Pr_{R_\beta \eta}$ for all $\beta \in \Lambda$. Suppose that $Pr_\xi \neq Pr_\eta$. It means that there exists a function $f(x, y) \in L^1(U)$ such that

$$Pr(\xi(x, y) = f(x, y)) \neq Pr(\eta(x, y) = f(x, y)),$$

and hence

$$|Pr(\xi(x, y) = f(x, y)) - Pr(\eta(x, y) = f(x, y))| = \delta > 0.$$

Denote by $f_1(x, y), f_2(x, y), \dots$ values of the random element $\xi(x, y)$, different from $f(x, y)$ (we assume that these values are numbered in the decreasing order of probabilities $Pr(\xi(x, y) = f_i(x, y))$), and by $g_1(x, y), g_2(x, y), \dots$ analogous values of $\eta(x, y)$ (which are also numbered in the decreasing order of probabilities $Pr(\eta(x, y) = g_i(x, y))$) (thus $f(x, y) \neq f_i(x, y)$ and $f(x, y) \neq g_i(x, y)$, $i = 1, 2, \dots$). For each fixed $i = 1, 2, \dots$ let A_i be the set of all $\beta \in \Lambda$, for which

$$R_\beta f(\gamma) \equiv R_\beta f_i(\gamma)$$

and, respectively, B_i be the set of all $\beta \in \Lambda$, for which

$$R_\beta f(\gamma) \equiv R_\beta g_i(\gamma).$$

Each of the sets A_i and B_i is finite (maybe empty). Indeed, if any set A_i is infinite then by lemma 1 we have $f(x, y) \equiv f_i(x, y)$. Similarly, if any set B_i is infinite, then $f(x, y) \equiv g_i(x, y)$.

Thus, the set

$$C_n = \bigcup_{i=1}^n (A_i \cup B_i)$$

is finite for any n , hence the set $\Lambda \setminus C_n$ is not empty. Consider arbitrary $\beta \in \Lambda \setminus C_n$. Since for this β we have $R_\beta f(\gamma) \neq R_\beta f_i(\gamma), i = 1, \dots, n$, then

$$\begin{aligned} Pr(\xi(x, y) = f(x, y)) &\leq Pr(R_\beta \xi(\gamma) = R_\beta f(\gamma)) \\ &\leq Pr(\xi(x, y) = f(x, y)) + \sum_{i=n+1}^{\infty} Pr(\xi(x, y) = f_i(x, y)), \end{aligned}$$

and so

$$\begin{aligned} |Pr(R_\beta \xi(\gamma) = R_\beta f(\gamma)) - Pr(\xi(x, y) = f(x, y))| \\ \leq \sum_{i=n+1}^{\infty} Pr(\xi(x, y) = f_i(x, y)) < \varepsilon_{1,n}, \end{aligned}$$

where $\varepsilon_{1,n} \rightarrow 0$ when $n \rightarrow \infty$.

Similarly,

$$\begin{aligned} & |Pr(R_\beta \eta(\gamma) = R_\beta f(\gamma)) - Pr(\eta(x, y) = f(x, y))| \\ & \leq \sum_{i=n+1}^{\infty} P(\eta(x, y) = g_i(x, y)) < \varepsilon_{2,n}. \end{aligned}$$

Since $Pr(R_\beta \xi(\gamma) = R_\beta f(\gamma)) = Pr(R_\beta \eta(\gamma) = R_\beta f(\gamma))$, then

$$\begin{aligned} & |Pr(\xi(x, y) = f(x, y)) - Pr(\eta(x, y) = f(x, y))| \\ & \leq |Pr(R_\beta \xi(\gamma) = R_\beta f(\gamma)) - Pr(\xi(x, y) = f(x, y))| \\ & \quad + |Pr(R_\beta \eta(\gamma) = R_\beta f(\gamma)) - Pr(\eta(x, y) = f(x, y))| \leq 2\varepsilon_n, \end{aligned}$$

where $\varepsilon_n = \max(\varepsilon_{1,n}, \varepsilon_{2,n})$ can be made arbitrarily small. But by our assumption

$$|Pr(\xi(x, y) = f(x, y)) - Pr(\eta(x, y) = f(x, y))| = \delta > 0.$$

This contradiction proves the theorem.

6. Classification of projections. Theorem 1 from the previous section shows that in class T it is possible to reconstruct distribution of bivariate random function from distributions of projections. In this section we develop an algorithm, which allows to sort registered projections by the groups corresponding to different states of a random function.

For the sake of simplicity we will consider functions from class T , having at most two different states. Algorithm can be directly generalized to any finite number of states, and for the case of denumerable number of states we can “truncate” distributions of projections as it is done in [4]. In addition, since in the frames of considered problem functions $f_i(x, y)$ describe distribution of object density, we can assume that they are nonnegative. We will also assume that these functions are normalized, i. e. integrals of them are equal to 1 (when these conditions are met, functions $f_i(x, y)$ are probability densities).

So, let $\xi(x, y)$ be a random function that takes values $f_1(x, y)$ and $f_2(x, y)$ with probabilities p_1 and p_2 . We assume that distributions of $R_\beta \xi(\gamma)$ are known for each $\beta \in \Lambda \subseteq [0, 2\pi)$. I. e. for each β we know functions $R_\beta f_i(\gamma)$, $i = 1, 2$, which are projections of functions $f_i(x, y)$, $i = 1, 2$, and are realized with probabilities p_1 and p_2 respectively. As a matter of fact, we don’t know in advance, which realization of projection corresponds to which realization of

a random function, i.e. it is possible that for some β function $R_\beta f_1(\gamma)$ is a projection of $f_2(x, y)$, and $R_\beta f_2(\gamma)$ is a projection of $f_1(x, y)$. We have to sort functions $R_\beta f_i(\gamma)$, $i = 1, 2$, for all $\beta \in \Lambda$ by the groups, so that each group contains realized projections that correspond to the same state of a random function.

If $p_1 \neq p_2$ then such classification can be performed according to probabilities of states of projections, i.e. for each $\beta \in \Lambda$ value $R_\beta f_i(\gamma)$, which is realized with probability p_1 , we ascribe to the first group, and value $R_\beta f_i(\gamma)$, which is realized with probability p_2 , we ascribe to the second group.

In case $p_1 = p_2 = \frac{1}{2}$, projection classification algorithm is based on the construction of approximations for functions $H^{(i)}(\beta, \varphi)$, $i = 1, 2$, which are defined for $f_1(x, y)$ and $f_2(x, y)$ by expression (2). If $f_1(x, y) \not\equiv f_2(x, y)$, then for almost every fixed φ , for which functions $H^{(i)}(\beta, \varphi)$ are defined by expression (4), they may coincide only for finite number of β . Hereinafter we will consider only those φ . Let $\Lambda = \bigcup_{k=1}^K \Lambda_k$, and let φ_k , $k = 1, \dots, K$, be such that for all $\beta \in \Lambda_k$ and some $\delta > 0$ we have $D \sin(\varphi_k - \beta) > 1 + \delta$. We assume that intersection of adjacent Λ_k is not empty and for $\beta \in \Lambda_k \cap \Lambda_{k+1}$ we have $|H^{(1)}(\beta, \varphi_k) - H^{(2)}(\beta, \varphi_k)| > A$ and $|H^{(1)}(\beta, \varphi_{k+1}) - H^{(2)}(\beta, \varphi_{k+1})| > A$ for some $A > 0$. We will perform classification on each Λ_k separately. Let

$$H_n^{(i)}(\beta, \varphi) = \sum_{m=0}^n (D \sin(\varphi - \beta))^{-(m+1)} p_m^{(i)}(\varphi), \quad i = 1, 2.$$

If all conditions on functions $f_i(x, y)$, described in the beginning of the section, are met, then for all m and φ we have $|p_m^{(i)}(\varphi)| \leq 1$. Thus for all $\beta \in \Lambda_k$ we have

$$|H^{(i)}(\beta, \varphi_k) - H_n^{(i)}(\beta, \varphi_k)| \leq \frac{1}{(1 + \delta)^n}. \quad (5)$$

Let $s_j = (D \sin(\varphi_k - \beta_j))^{-1}$, and let β_j , $j = 0, \dots, n$, be the collection of Tchebyshev interpolation points for Λ_k (see [7]). Using expression (2), we can calculate values of $H^{(i)}(\beta, \varphi_k)$ in points β_j . There are 2^{n+1} ways to distribute values $H^{(i)}(\beta_j, \varphi_k)$, $i = 1, 2$, among two groups. Let us denote by P the set of all such distribu-

tions and solve linear systems

$$\sum_{m=0}^n s_j^m a_m^{(h)} = H^{(h)}(\beta_j, \varphi_k), \quad j = 0, \dots, n$$

for every possible distribution h from P . Following ideas of [3], we can obtain 2^{n+1} functions of the form

$$G_n^{(h)}(\beta, \varphi_k) = \sum_{l=0}^n \frac{\prod_{j \neq l} ((\sin(\varphi_k - \beta))^{-1} - (\sin(\varphi_k - \beta_j))^{-1})}{\prod_{j \neq l} ((\sin(\varphi_k - \beta_l))^{-1} - (\sin(\varphi_k - \beta_j))^{-1})} H^{(h)}(\beta_j, \varphi_k) \quad (6)$$

(here superscript h denotes chosen distribution of values $H^{(i)}(\beta_j, \varphi_k)$, $i = 1, 2$, among two groups), from which we must choose approximations for functions $H^{(i)}(\beta, \varphi_k)$, $i = 1, 2$. To evaluate inaccuracy of approximations for functions $H^{(i)}(\beta, \varphi_k)$, $i = 1, 2$, we use well-known estimate of Tchebyshev interpolation (see [7]). From (5) it follows that for function $G_n^{(h)}(\beta, \varphi_k)$, which approximates $H_n^{(i)}(\beta, \varphi_k)$ for $i = 1$ or $i = 2$, we must have

$$|G_n^{(h)}(\beta, \varphi_k) - H_n^{(1)}(\beta, \varphi_k)| \leq \frac{1}{(1+\delta)^n} \left(8 + \frac{4}{\pi} \ln(n+1) \right)$$

or

$$|G_n^{(h)}(\beta, \varphi_k) - H_n^{(2)}(\beta, \varphi_k)| \leq \frac{1}{(1+\delta)^n} \left(8 + \frac{4}{\pi} \ln(n+1) \right)$$

for all $\beta \in \Lambda_k$. Hence for all $\beta \in \Lambda_k$ we must have

$$|G_n^{(h)}(\beta, \varphi_k) - H^{(1)}(\beta, \varphi_k)| \leq \frac{1}{(1+\delta)^n} \left(9 + \frac{4}{\pi} \ln(n+1) \right) \quad (7')$$

or

$$|G_n^{(h)}(\beta, \varphi_k) - H^{(2)}(\beta, \varphi_k)| \leq \frac{1}{(1+\delta)^n} \left(9 + \frac{4}{\pi} \ln(n+1) \right). \quad (7'')$$

Starting with some n there are only 2 distributions h_1 and h_2 from P , for which these inequalities hold for all $\beta \in \Lambda_k$. For those h_1 and h_2 functions $G_n^{(h_i)}(\beta, \varphi_k)$, $i = 1, 2$, are taken to be approximations for $H^{(i)}(\beta, \varphi_k)$.

Then, calculating values $H^{(i)}(\beta, \varphi_k)$, $i = 1, 2$, with the use of (2) for each $\beta \in \Lambda_k$, we ascribe projections to the first or second

group, depending on the closeness of these values to the values of functions $G_n^{(h_i)}(\beta, \varphi_k)$, $i = 1, 2$, in point β .

Since intersection of adjacent Λ_k is not empty, for each $\beta \in \Lambda_k \cap \Lambda_{k+1}$ values of $H^{(i)}(\beta, \varphi)$, $i = 1, 2$, for $\varphi = \varphi_k$ and $\varphi = \varphi_{k+1}$ will be grouped correctly by the states of random function, because they are calculated for the same realization of projection and for some A we have $|H^{(1)}(\beta, \varphi_k) - H^{(2)}(\beta, \varphi_k)| > A$ and $|H^{(1)}(\beta, \varphi_{k+1}) - H^{(2)}(\beta, \varphi_{k+1})| > A$. (Remark: this condition is required because algorithm of classification may fail in the neighborhood of intersection points of functions $H^{(1)}(\beta, \varphi_k)$ and $H^{(2)}(\beta, \varphi_k)$. As n increases, the total length of these neighborhoods tends to zero.) Thus the “linkage” of groups for adjacent Λ_k is implemented, and classification is performed on the whole set Λ .

When all projections are grouped we can reconstruct each state of random function (and hence its distribution), using conventional techniques of computer tomography.

Described algorithm is exact in the sense that if projections are given exactly, then it is possible, in principle, to approximate functions $H^{(i)}(\beta, \varphi_k)$ arbitrarily close and use their values to group projections. In practice, however, it is impossible because of at least two reasons. First reason lies in the basis of method itself, because when expanding $H^{(i)}(\beta, \varphi_k)$ into series, we keep only finite number of terms. Second reason is associated with inability of exact registration of projections. If projections are given with some error, then as n increases, estimates in the right-hand sides of (7') and (7'') will decrease only to the certain limit, and after that they will begin to increase. If projections $R_\beta(\gamma)$ are given with error, whose level does not exceed ε , then, using (3) and the fact that if error of $R_\beta(\gamma)$ does not exceed ε , then error of $P_\varphi(t)$ also does not exceed ε , we can obtain the following estimates for approximations of $H^{(i)}(\beta, \varphi_k)$:

$$|G_n^{(h_i)}(\beta, \varphi_k) - H^{(i)}(\beta, \varphi_k)| \leq \left(\frac{2\varepsilon}{\delta} + \frac{1}{(1+\delta)^n} \right) \left(9 + \frac{4}{\pi} \ln(n+1) \right). \quad (8)$$

References

1. Liu W., Frank J. Estimation of variance distribution in three-dimensional reconstruction. I. Theory // J. Opt. Soc. Am. A. 1995. V. 12. P. 2615–2627.

2. *Ushakov V. G., Ushakov N. G.* Reconstruction of probabilistic characteristics of multivariate random functions from their projections // Vestn. Mosk. Univ. Ser. 15: Comput. Math. Cybern. 2001. No. 4. P. 32–39 (in Russian).
3. *Shestakov O. V.* On the uniqueness of reconstruction of probabilistic characteristics of multivariate random functions from probabilistic characteristics of their projections // Vestn. Mosk. Univ., Ser. 15: Comput. Math. Cybern. 2003. No. 3. P. 37–41 (in Russian).
4. *Shestakov O. V.* An algorithm to reconstruct probabilistic distributions of multivariate random functions from the distributions of their projections // Journal Of Mathematical Sciences. 2002. V. 112, No. 2. P. 4198–4204.
5. *Troitskiy I. N.* Statistical theory of tomography. — Moscow, Radio i Svyaz, 1989 (in Russian).
6. *Hamaker C., Smith K. T., Solmon D. C., Wagner S. L.* The divergent beam X-ray transform // Rocky mountain journal of mathematics. 1980. V. 10, No. 1. P. 253–283.
7. *Goncharov V. L.* Theory of Interpolation and Approximation of Functions. — Moscow, 1933 (in Russian).

ESTIMATION OF DELAY DISTRIBUTION IN HIV DYNAMICS

Anastasia Ushakova

Department of Mathematical Sciences, Norwegian
University of Science and Technology, Trondheim,
Norway

In this paper we present two methods of estimating a delay distribution in biological dynamical systems. The model of HIV infection serves as an example of such systems. The first method is based on parametric approach and on approximation of the delay density by a gamma-density. The second method is nonparametric and is based on solution of a convolution equation with selection of the regularization parameter from a parametric start.

1. Introduction

One of the main problems in HIV infection research is understanding the events that occur during the long asymptomatic period of the disease. This period, which usually starts after a few weeks (or in some cases months) after infection, can last for many years and inevitably turns into the AIDS stage, when the viral load increases rapidly and the CD+ T-cell count declines to a point at which the immune system fails to provide protection against opportunistic infections. To describe quantitatively the processes under consideration, one needs mathematical models which, on the one hand, have to be adequate and, on the other hand, should be simple enough. There are a number of such models, see for example [1].

Let $T(t)$, $I(t)$, V_I и V_{NI} be respectively the density of non-infected target cells, the density of productively infected cells, the concentration of infections virus and the concentration of virus that has been rendered non-infections by the protease inhibitors. One of the most used and relatively simple models, which was developed in [2], is as follows. First it is supposed that $T(t)$ is constant $T(t) = T$. Then

$$\frac{dI}{dt} = kTV_I(t) - \delta I(t), \quad (1)$$

$$\begin{aligned}\frac{dV_I}{dt} &= (1 - \eta)pI(t) - cV_I(t), \\ \frac{dV_{NI}}{dt} &= \eta pI(t) - cV_{NI}(t),\end{aligned}$$

where k is the infection rate constant, p the rate at which a productively infected cell release virions, and η the drug efficacy (if $\eta = 1$, the drug is assumed to be absolutely effective so that all virions produced after drug takes effect are noninfectious). Productively infected cells die at a rate per cell δ and plasma virions are cleared at a rate c per virion. In this model however it is not taken into account that there is a certain time lapse between viral entry into a target cell and the production of new virus particles — an intracellular delay. A number of authors, see in particular [3]–[5], demonstrated that this can misrepresent the real process, and therefore the intracellular delay has to be taken into consideration. On the other hand, conversion of a newly infected cell into a productively infected cell is a complicated multi-step process that can last from one to several days, and therefore, fixed delays are not realistic. A model containing a distributed intracellular delay was suggested in [4]. According to this model equation (1) is replaced by the following equation

$$\frac{dI}{dt} = kT \int_0^{\infty} V_I(t - t')f(t') dt' - \delta I(t). \quad (2)$$

Here the function $f(t)$ — the density of the delay distribution — has to be estimated from observations.

2. Parametric estimation. Gamma model

Equation (2) is a special case of convolution equations of the first kind. Since methods developed in this work do not use specifics of equation (2), we consider it in the general form, namely

$$\int_0^{\infty} K(t - s)z(s)ds = u(t), \quad (3)$$

where $K(t)$ and $u(t)$ are observed functions and $z(t)$ is estimated function (in our case — the density of the delay distribution). Functions $K(t)$ and $u(t)$ are supposed to be integrable, non-negative and that they tend to zero as $t \rightarrow \pm\infty$.

For biological objects, a gamma-distribution is often a good approximation for the delay distribution [6]. In this section we suppose that $z(t)$ is a gamma-density i. e. it has form

$$z(t) = \frac{\beta^\alpha}{\Gamma(\alpha)} t^{\alpha-1} e^{-\beta t}, \quad t > 0; \quad \alpha > 0, \beta > 0.$$

In this case the problem is reduced to estimation of two parameters, α and β . Suppose that functions $K(t)$ and $u(t)$ are observed with random errors at time points t_1, \dots, t_n .

It was suggested in [4] to estimate the parameters using the least square method. This, however, leads to minimization (with respect to α and β) of the expression

$$S(\alpha, \beta) = \sum_{j=1}^n \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^\infty K(t_j - s) s^{\alpha-1} e^{-\beta s} ds - u(t_j) \right)^2,$$

and therefore one needs to solve strongly nonlinear equations, very sensitive to measurement errors and to the replacement of the continuous model by the discrete one. Here we present a method which on the one hand is very simple and on the other hand is quite stable with respect to measurement errors.

Let μ_z and σ_z^2 be the expectation and the variance of the probability density $z(t)$. We have

$$\alpha = \frac{\mu_z^2}{\sigma_z^2}, \quad \beta = \frac{\mu_z}{\sigma_z^2}.$$

Denote

$$\mu_K = \frac{\int_0^\infty t K(t) dt}{\int_0^\infty K(t) dt}, \quad \mu_u = \frac{\int_0^\infty t u(t) dt}{\int_0^\infty u(t) dt} \quad (4)$$

and

$$\sigma_K^2 = \frac{\int_0^\infty (t - \mu_K)^2 K(t) dt}{\int_0^\infty K(t) dt}, \quad \sigma_u^2 = \frac{\int_0^\infty (t - \mu_u)^2 u(t) dt}{\int_0^\infty u(t) dt}. \quad (5)$$

Then

$$\mu_z = \mu_u - \mu_K, \quad \sigma_z^2 = \sigma_u^2 - \sigma_K^2$$

and hence

$$\alpha = \frac{(\mu_u - \mu_K)^2}{(\sigma_u^2 - \sigma_K^2)}, \quad \beta = \frac{\mu_u - \mu_K}{(\sigma_u^2 - \sigma_K^2)}. \quad (6)$$

Suppose that the points t_1, \dots, t_n form a uniform grid with the step size h , and that this grid covers a region such that one may neglect functions $K(t)$ and $u(t)$ outside this region. One can estimate parameters μ_K , μ_u , σ_K^2 and σ_u^2 replacing integrals in (4) and (5) by the corresponding integral sums, i. e. as follows

$$\hat{\mu}_K = \frac{\sum_{i=1}^n t_i K(t_i)}{\sum_{i=1}^n K(t_i)}, \quad \hat{\mu}_u = \frac{\sum_{i=1}^n t_i u(t_i)}{\sum_{i=1}^n u(t_i)}$$

and

$$\hat{\sigma}_K^2 = \frac{\sum_{i=1}^n (t_i - \hat{\mu}_K)^2 K(t_i)}{\sum_{i=1}^n K(t_i)}, \quad \hat{\sigma}_u^2 = \frac{\sum_{i=1}^n (t_i - \hat{\mu}_u)^2 u(t_i)}{\sum_{i=1}^n u(t_i)}.$$

Substitution of these estimators in (6) then yields

$$\hat{\alpha} = \frac{(\hat{\mu}_u - \hat{\mu}_K)^2}{(\hat{\sigma}_u^2 - \hat{\sigma}_K^2)}, \quad \hat{\beta} = \frac{\hat{\mu}_u - \hat{\mu}_K}{(\hat{\sigma}_u^2 - \hat{\sigma}_K^2)}.$$

Estimators $\hat{\alpha}$ and $\hat{\beta}$ display quite a good performance. In Table 1 we present estimates of the standard error of estimators $\hat{\alpha}$ and $\hat{\beta}$ obtained by simulation. It is supposed that

$$K(t_j) = K_0(t_j) + \xi_j, \quad u(t_j) = u_0(t_j) + \eta_j$$

where $K_0(t)$ and $u_0(t)$ are the exact kernel and the exact right hand side of equation (3), $\xi_1, \dots, \xi_n, \eta_1, \dots, \eta_n$ are random errors. Random variables ξ_1, \dots, ξ_n as well as η_1, \dots, η_n are independent and identically distributed. ξ_i and η_i have normal distributions with zero mean and variances σ_ξ^2 and σ_η^2 respectively. Variances of $\hat{\alpha}$ and $\hat{\beta}$ are denoted by σ_α^2 and σ_β^2 . The kernel $K(t)$ is normal

$$K(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(t-4)^2}{2}}$$

Table 1

Performance of the estimators

α	β	$\sigma_\xi/\max K(t)$	$\sigma_\eta/\max u(t)$	$\hat{\sigma}_\alpha/\alpha$	$\hat{\sigma}_\beta/\beta$
1	1	0.025	0.032	0.080	0.067
1	1	0.050	0.064	0.171	0.143
1	1	0.100	0.128	0.371	0.299
2	1	0.025	0.040	0.033	0.031
2	1	0.050	0.080	0.063	0.061
2	1	0.100	0.155	0.131	0.125
2	2	0.025	0.030	0.097	0.087
2	2	0.050	0.060	0.225	0.202
2	2	0.100	0.119	0.533	0.468
3	1	0.025	0.045	0.009	0.009
3	1	0.050	0.091	0.020	0.021
3	1	0.100	0.182	0.042	0.046
3	1	0.251	0.456	0.106	0.110
3	2	0.025	0.032	0.085	0.078
3	2	0.050	0.064	0.180	0.165
3	2	0.100	0.128	0.396	0.355

(this distribution can be considered as concentrated on the positive half-line).

Note that for inverse problems of this kind, the accuracy $\sim \sqrt{\varepsilon}$ (as $\varepsilon \rightarrow 0$) of the solution, where ε is the accuracy of the data, is usually considered as a good result. Here the result is much better.

3. Nonparametric estimation of delay distribution

Parametric approach is very convenient for estimation of the delay distribution but, using it, one can lose some important features of the density to be estimated. For example, all gamma densities are unimodal. While presence of two or more local maxima is, on the one hand, quite possible and, on the other hand, can reflect essential peculiarities of the process under consideration. In this

section, we briefly consider nonparametric estimation of the delay distribution.

Equation (3) is an integral equation of the first kind of the convolution type. Methods of solution of such equations are well developed and use regularization, since the problem is ill-posed. Here we use the standard technique. However, the main part of the solution — selection of the regularization parameter — can be worked out using specifics of the applied problem under consideration. In our case, parametric approximation can cause a loss of important features of the density to be estimated, but it allows one to estimate the measurement error or smoothness level of the estimated density quite well.

Denote Fourier transforms of functions $K(t)$, $z(t)$ and $u(t)$ by the same letters but of the variable ω . Regularized solution has form (see for example [7])

$$z_\delta(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{K(-\omega)u(\omega)e^{-i\omega t}}{|K(\omega)|^2 + \delta M(\omega)} d\omega,$$

where $M(\omega)$ is an even nonnegative function, such that $M(0) \geq 0$, $M(\omega) > 0$ for $\omega \neq 0$ and $M(\omega) \geq c > 0$ for sufficiently large $|\omega|$, satisfying some regularity conditions (see [7] for details), and δ is the regularization parameter (a positive number). One can put for example

$$M(\omega) = 1 - e^{-c\omega^2},$$

where c is chosen from the condition that functions $\exp(-c\omega^2)$ and $K(\omega)$ have approximately the same range.

If the deviation γ of the observed right hand side of equation (3) $u(t)$ from the exact one $u_T(t)$ (in some metric $\rho(\cdot, \cdot)$) was known: $\rho(u, u_T) = \gamma$, one could select the regularization parameter from the following condition

$$\rho(K^* z_\delta, u) = \gamma.$$

This deviation is unknown, but can be estimated from the data. We suggest to select the regularization parameter as follows. On the first stage we use parametric approach and approximate $z(t)$ by a gamma-density. Parameters of the gamma-density α and β are estimated using the technique described in the previous section.

Let $\widehat{z}(t)$ be the obtained estimate. Set

$$\widehat{u}(t) = \int_0^{\infty} K(t-s)\widehat{z}(s) ds.$$

The measurement error (in L^2 metrics) of the right hand side of equation (3) is now estimated by

$$\Delta = \int_0^{\infty} (\widehat{u}(t) - u(t))^2 dt.$$

The regularization parameter δ is selected as solution of the equation

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{\delta^2 M^2(\omega) |u(\omega)|^2}{(|K(\omega)| + \delta M(\omega))^2} d\omega = \Delta,$$

because the left hand side of this equation is

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} |K(\omega)z_{\delta}(\omega) - u(\omega)|^2 d\omega = \int_0^{\infty} (K^*z(t) - u(t))^2 dt.$$

Another way to use a parametric start in nonparametric estimation of the delay distribution consists in estimation of some functional of smoothness of the estimated density, for example the total variation. This functional is estimated on the basis of parametric approximation as it was described above.

References

1. *Nowak M. A., May R. M.* Virus dynamics: mathematical principles of immunology and virology. — Oxford University Press, 2000.
2. *Perelson A. S., Neumann A. U., Markowitz M., Leonard J. M., Ho D. D.* HIV-1 dynamics in vivo: Virion clearance rate, infected cell life-span, and viral generation time // Science. 1996. V. 271. P. 1582–...
3. *Herz A. V. M., Bonhoeffer S., Anderson R. M., May R. M., Nowak M. A.* Viral dynamics in vivo: Limitations on estimates of intracellular delay and virus decay // Proc. Natl. Acad. Sci. USA. 1996. V. 93. P. 7247–...

4. *Mittler J. E., Bernhard S., Neumann A. U., Perelson A. S.* Influence of delayed viral production on viral dynamics in HIV-1 infected patients // *Mathematical Biosciences*. 1998. V. 152. P. 143–163.
5. *Nelson P. W., Perelson A. S.* Mathematical analysis of delay differential equation models of HIV-1 infection // *Mathematical Biosciences*. 1998. V. 179. P. 73–94.
6. *MacDonald N.* Biological delay systems: linear stability theory. — Cambridge University, Cambridge, 1989.
7. *Tikhonov A. N., Arsenin V. Y.* Solutions of ill-posed problems. — N. Y., Wiley, 1977.

ON LINGUISTIC CLASSIFICATION OF BACTERIAL GENOMES

Zeev Volkovich

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

Valery Kirzhner

Institute of Evolution, University of Haifa, Haifa, Israel

Zeev Barzily

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

The paper is devoted to classification of 185 full prokaryote genomes using a modification of the compositional spectra method. This modification suggests separate calculation of the compositional spectra for coding and non-coding subsequences of the genome. For each subsequence, the corresponding vector, in Euclidian space, can be obtained using certain manipulations of the compositional spectra. This allows analyzing the structure of genome and determining the most probable number of genome clusters without any additional information. Our clustering method is based on the application of the external indexes of partitions agreement and the number of the misclassified items within repeated partitions. A biological justification, for the four and the two letters alphabets, substantiates the appropriateness of the outcomes acquired.

1. Background

With the increasing number of full genome sequences available, new challenges are emerging in the field of computational biology. The existence of a long single contig, representing the full genomic content of a bacterial strain, allows the study of a genome as a whole and not as a collection of genes. From this point of view, the question is not to analyze particular features of a protein, or a family of proteins, but to consider the global properties of the genetic text of a bacterium. The first important step of a global genome analysis is the description of general rules that allow the merging of the different types of information within the physiological and environmental contexts ([1] Rocha et al., 1988). In this matter, it is interesting to compare between two of the main functional genomic

structures -coding and non-coding parts of full genome sequences. In prokaryotes, non-coding regions and genes seem to have evolved in different regions ([2] Rogozin et al., 2002). The dynamic organization of non-coding DNA suggests a feedback loop which can influence codon usage and can stabilize the chromosome's chromatin pattern ([3] Holmquist, 1989). Hierarchical selection theories show how selection can act on non-coding DNA, at the genome level, creating positional constrained DNA and contributing minimal genetic load at the individual level [3–6]. Comparing genome sequences, in particular discovering distances between sequences, is an important issue in bioinformatics. Different formalisms have been introduced to construct genome distances on the basis of various aspects of the genome (see, [7]). Many approaches for genomes representation by means of partial word frequencies can be found in the literature. For example, the identification of a specific site [8] and the genome annotation [9]. Article [10] determines, for each bacterium, a set of oligonucleotides uniquely characterizing the organism. A method for dendrograms construction using the word based approaches is implemented in [11, 12].

One of the approaches for representing the structure of a genome is termed Compositional Spectra (CS) [13, 14]. One of the benefits of this approach is that it is capable to construct an embedding of the genome sequences into a Euclidian space taking into account the genome textual structure. Particularly, it makes possible to estimate distances between whole genomes [15, 18] and to provide classifications of the genomes via the conventional clustering approaches (see, for example [15]).

In this paper, we apply a modification of the CS method, in which spectra for coding and non-coding parts of the whole genome are calculated separately. For each subsequence, the corresponding vector in Euclidian space can be obtained using certain manipulations with compositional spectra. This allows to analyze the structure of genome and determine the most probable number of genome clusters without any additional information. Our method for determining the true number of clusters is based on an application of the external indices of partitions agreement and the number of the misclassified items within repetitive clustering of the same dataset. Note that such methodology has been used in different versions (see, for example [16, 17]). However, we would like to emphasize that, to our knowledge, there has not been any similar effort made

before, in such a context, where the clustering is completed, by means of two different embeddings into geometrical spaces.

2. Methods

2.1. Compositional spectrum. In previous works [13, 14] the notion of *Compositional Spectrum* (CS) has been introduced as follows. Consider a word w of length L in the alphabet $\{A, T, C, G\}$ and sequence S in the same alphabet. Let sequence S contain a word x , which differs from word w no more than in r positions, hence, word x is an imperfect occurrence of w in S . This approximate matching is referred as “ r -mismatching”.

Let us look at a set $W = \{w_i\}$, $i = 1, \dots, n$, of n words of a length L and compare its elements with a given sequence S . Denote by m_i the number of r -imperfect occurrences produced by the word w_i over the sequence. The set $F(W, S)$, composed by the relative frequencies

$$f_i = \frac{m_i}{M}, \quad i = 1, \dots, n, \quad M = \sum_{i=1}^n m_i,$$

is named the “compositional spectrum” of the sequence S in consistence with the set W . In CS-analysis, we consider sufficiently long sequences S ($\approx 10^3$ – 10^5 bp) and usually use the values: $n = 200$ and $r = 2$. In the case of 4 letters alphabet we choose $L = 10$. With the two-letter alphabet ($A = C$, $T = G$; purine/pyrimidine, or R/Y), we take $L = 20$, so that the number of possible words will be the same, as with the 4-letter alphabet for $L = 10$. To produce a set W we employ a random number generator assuming equal probabilities of appearance of each of the four (or two) symbols at any current position. A distance between two sequences S and S' is now being constructed as a distance between their spectra $F(W, S)$ and $F(W, S')$, which can be derived in many ways. It has been shown [13, 14, 18, 19] that, from the biological point of view, the most adequate distances between species, produce correlation type functionals. The best tested functional, calculated for any two CS in [19], was $d = 1 - \rho$, where ρ is the Spearman rank correlation coefficient [20]. However, for this non metrical functional, it is complicated to use clustering methods for classifying of points in a metrical vector space. To overcome this difficulty, we suggest to calculate for each genome, the set of its distances d to all the

other genomes. Obviously, such a set can be viewed as a point in the corresponding metrical space. Hence, standard Euclidian type metrics describe, quite adequately, the distances between any pair of species. In what follows such metrics is called integrative metrics. It enables the applications of flexible partitioning procedures which are based on the Euclidian distance. The validity of the techniques described above is a consequence of the result, obtained in ([13, 14], because the distance between a pair of sequences, does not depend on the selection of W .

2.2. Clustering. Organization of data into a set of similar groups is an accepted and fundamental tool for data analysis. In data mining, clustering is one of the two most widespread techniques (the other one is classification) that are capable of providing a categorization of a set of items. Clustering produces a partition, of the data set, so that objects inside a cluster are similar to one another, but are unlike objects in other clusters. A similarity determines the group membership via a distance-like function which evaluates the resembling between two data points. Most widespread iterative clustering algorithms as k -means, k -medoids, EM and PAM are carried out, as a rule, in three steps. An initialization step is intended to set an initial partition. In the second step the data is partitioned to “best” possible clusters. It is made by assigning elements to clusters so that an objective function is optimized. The third step compares the current partition to the previous one. If the difference is less than the stopping parameter then the algorithm ends else it returns to the second step. The partitioning phase assigns a label to each item. This label identifies the cluster to which it belongs. Generally, labels’ values do not have specific meanings and can be permuted from one instance to another. One of input parameters of a clustering iterative routine is the suggested number of clusters in the considered dataset. Estimation of this number represents an ill-posed problem of crucial relevance in cluster analysis [21, 22]. For instance, the “correct” number of clusters in a dataset can depend on the scale in which the data is measured (see, for example, [23]). Generally speaking, approaches to this problem apply two methodologies. The first is based on geometrical features of partitions like within and between cluster dispersions, and the second is base on the cluster stability properties. For example, the Gap statistics method proposed in

[24] can be cited as one of the approaches belonging to the first methodology.

In our model the distances between the items are calculated resting up different parts of a genome with the appropriate embeddings into Euclidian spaces. Thus, a comprehensible geometrical interpretation of the objects is hardly attained. However, the duplicate clustering of each genome, by its distinct portions, yields the model's stability which is reflected by the consistency within repeated clustering procedures.

External indices of partitions agreement. External indices of partitions agreement yield a common tool for cluster stability assessing. In this contest the Rand coefficient [25], the Jaccard coefficient [21] and the Fowlkes and Mallows coefficient [26] can be mentioned. The Clest approach [16] uses such coefficients for the cluster validation determination. A method offered in [27] characterizes a cluster stability event by the distribution of pair-wise agreements between partitions, built by means of the external indexes on sub-samples of the data. In this section we state several facts about external indices of partitions agreement.

Let $X = \{x_1, \dots, x_n\}$ be a finite set. Consider a partition $\Pi_k = \{\pi_1, \dots, \pi_k\}$ of the set, i.e.,

$$\pi_1 \cup \dots \cup \pi_k = X \quad \text{and} \quad \pi_i \cap \pi_j = \emptyset \quad \text{if} \quad i \neq j.$$

The elements of the partition are named as clusters. The calculation of the external indexes is based on cross-tabulation, or contingency tables, which are composed by the co-occurrences of objects belonging to clusters in partitions $\Pi_r^{(1)}$ and $\Pi_c^{(2)}$ of X . Two partitions are identical if and only if every cluster in $\Pi_r^{(1)}$ is also a cluster in $\Pi_c^{(2)}$. Namely, the clusters in the partitions may only be differently labeled in the clusters' designates. Let us consider the following labeling representation of a given partition:

$$C = \{c_{ij}\} = \begin{cases} 1 & \text{if } x_i \text{ and } x_j \text{ belong to the same cluster } x_i \neq x_j, \\ 0 & \text{otherwise;} \end{cases}$$

$$i, j = 1, \dots, n.$$

Let C_1 and C_2 be two binary matrices of this kind. For $k, l = 0, 1$ we introduce the quantities N_{kl} of the number of items' pairs which take value k in C_1 and value l in C_2 respectively. For instance, N_{11}

is the number of pairs belonging to the same cluster in $\Pi_r^{(1)}$ and to the same cluster in $\Pi_c^{(2)}$. The following coefficients are often used:

- $Rand\left(\Pi_r^{(1)}, \Pi_c^{(2)}\right) = \frac{N_{00} + N_{11}}{N_{00} + N_{10} + N_{01} + N_{11}}$ (Rand coefficient);
- $JD\left(\Pi_r^{(1)}, \Pi_c^{(2)}\right) = \frac{N_{11}}{N_{01} + N_{10} + N_{11}}$ (Jaccard or Jain and Dubes coefficient);
- $FM\left(\Pi_r^{(1)}, \Pi_c^{(2)}\right) = \frac{N_{11}}{\sqrt{(N_{11} + N_{10})(N_{11} + N_{01})}}$ (Fowlkes-Mallows coefficient).

We ignore here the trivial cases where the denominators equal to zero. Due to the fact N_{kl} are nonnegative, all introduced indices receive values in the interval $[0, 1]$. Two partitions coincide if and only if $N_{10} = N_{01} = 0$. When partitions are equal, and $N_{11} \neq 0$, then all indices achieve their maximal value 1. An external index is frequently standardized in such a way that its expected value is 0, if the partitions are random, and 1 when they correspond perfectly. The general formula is offered to standardize an index is:

$$Ind' = \frac{(Ind - E(Ind))}{(Ind_{\max} - E(Ind))}.$$

Particularly, the adjusted Rand [28] index equals:

$$Rand' = \frac{N_{11} - (N_{11} + N_{10})(N_{11} + N_{01})/N_{00}}{(N_{11} + N_{10} + N_{01})/2 - (N_{11} + N_{10})(N_{11} + N_{01})/N_{00}}. \quad (1)$$

The most popular null model assumes that the mentioned contingency table is created from the generalized hyper-geometric distribution and that the two partitions are mutually independent. In this case the adjusted index has to be zero. Values close to zero specify that from each partition, nothing can be forecasted about the other.

The number of the misclassified items. Another similarity measure, between two partitions, has been offered in [17, 29]. This measure can be defined as the quantity of the misclassified items, between the two partitions $\Pi_k^{(1)}$ and $\Pi_k^{(2)}$, of the same set X . An element is considered misclassified if it belongs to different clusters in $\Pi_k^{(1)}$ and $\Pi_k^{(2)}$. Such elements can be described in terms of the labeling functions α_1 and α_2 , from X to $C_k = [1, \dots, k]$,

defined as $\alpha_i(x) = c$, if and only if $x \in \pi_c$, $i = 1, 2$, $c = 1, \dots, k$. One of the problems, which usually arise here, is induced by the inherent permutation symmetry of clustering algorithms. Explicitly, the cluster labels are arbitrarily permuted. A matching between the labels can be found by resolving the task

$$D_k(\alpha_1, \alpha_2) = \min_{\psi \in \Psi^*} \sum_{x \in X} \chi(\alpha_1(x) \neq \psi(\alpha_2(x))), \quad (2)$$

where Ψ^* is the set of all possible permutations of C_k and χ is an indicator function of the event $\alpha_1(x) \neq \psi(\alpha_2(x))$. Note, we do not need to test all of $k!$ possible permutation, because this problem can be expressed as a partial case of the minimum weighed perfect bivariate matching problem. Computational complexity for solving of this problem by the well known Hungarian method (see, [30]) is $O(k^3)$. The range of values of $D_k(\alpha_1, \alpha_2)$ depends on k . Consequently, in order to compare $D_k(\alpha_1, \alpha_2)$, for different values of k , normalizing must be applied. Here, it is performed by

$$\tilde{D}_k(\alpha_1, \alpha_2) = \frac{D_k(\alpha_1, \alpha_2)}{E(D_k(\rho_1, \rho_2))}, \quad (3)$$

where ρ_1 and ρ_2 are random independent labelings uniformly distributed over C_k . Normalizing yields values of $D_k(\alpha_1, \alpha_2)$ independent of k . The expectation in the denominator can be calculated by simulation.

Algorithm description. Let us suppose that there are two embedding functions, E_1 and E_2 of the set X into the Euclidian space R^m . Assume that a clustering algorithm Cl is available. Among the input parameters of the algorithm a clustered set and a supposed number of clusters are included. The output is the items labels. Our algorithm can be described in the following form:

1. for $k = 2$ to k^*
2. for $t = 1$ to T
3. $\alpha_1(x) = Cl(E_1(X), k)$
4. $\alpha_2(x) = Cl(E_2(X), k)$
5. Calculate the value of an adjusted external index $Ind'_k(t)$ according to α_1 and α_2
6. Calculate the normalized quantity of misclassified items $\tilde{D}_k(\alpha_1, \alpha_2)(t)$ according to (2) and (3)
7. end for

8. $Ind'_k = \text{mean}(Ind'_k)$
9. $D_k = \text{mean}(\tilde{D}_k(\alpha_1, \alpha_2))$
10. end for
11. The true number of clusters is chosen according to the value of k yielding the maximum of the index Ind'_k or the minimum of D_k .

Here:

- k^* is a predefined parameter, which designates the maximal number of clusters to be tested;
- T is the iteration number of the indexes averaging.

It is well known, that the outcomes of an iterative clustering algorithm are strongly dependent on an initial partition. In order to avoid this difficulty, we average the indexes values. Note, that it is expected that the two criteria should yield the same value of k .

Clustering outcomes. We demonstrate the performance of our approach by analyzing a dataset consisting of 185 complete prokaryotic genomes. For the purposes of the study, every genome is transformed into two pseudo-sequences. The first one presents all coding sequences (*CDS*) and has been obtained by conjunction with the *CDS* preserving the order of genes and converting complementary genes to the correct “sense” format. The second pseudo-sequence *Sn* was obtained by conjugation of all non-coding fragments. In order to avoid artifacts, nonsense symbols out of the alphabet $\{A, T, C, G\}$ were inserted to separate gene fragments. Such insertions prevent considering non-existing words in the overall calculation. Furthermore, every genome is transformed into two CS yielding four possible distances $D_4(CDS)$ and $D_4(Sn)$, for the 4-letter alphabet, and $D_2(CDS)$ and $D_2(Sn)$ for the 2-letter purine/pyrimidine alphabet. Thus we get for the four-letter alphabet and a random set W_4 two spectra $F(W_4, CDS)$ and $F(W_4, Sn)$. Similarly, for the two-letter alphabet, and a set of two-letter words W_2 , two spectra $F(W_2, CDS)$ and $F(W_2, Sn)$ are defined. Next, the distances between spectra built on the same alphabet and the same genome’s part are calculated via the Spearman correlations, as it mentioned earlier. Hence, four distances matrices $D4(CDS)$ and $D4(Sn)$, for the 4-letter alphabet, and $D2(CDS)$ and $D2(Sn)$, for the 2-letter purine/pyrimidine alphabet, are resulted.

We represent each genome four times as a vector having 185 dimensions which are the distances to each of the genomes. We fur-

ther consider two pairs of embeddings, produced by the coding and the non-coding parts, for the 4-letter alphabet and for the 2-letter alphabet, respectively. As it was mentioned above, such an approach enables to obtain repeated clusterizations of the genome set.

For the clusterization we use the regular k -means algorithm based on the squared Euclidian distance. The Rand adjusted external index is employed for the partitions' similarity characterization. An application of our algorithm for the determination of "true" number of clusters, with the parameters $k^* = 9$ and $T = 200$, clearly indicates the presence of three clusters in all considered datasets. Typical graphs of the Rand adjusted external index and the normalized quantity of the misclassified items are presented below.

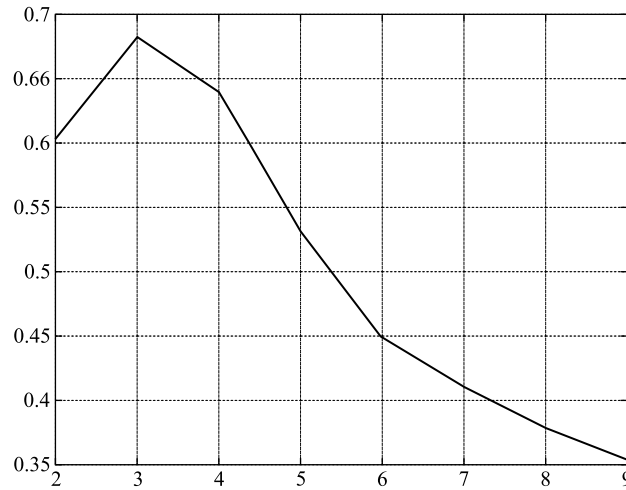


Fig. 1. Graph of the adjusted Rand coefficient

The final partitions for 3 clusters have been found by means of the regular k -means algorithm applying an additional initialization procedure offered in [19]. The initialization procedure eliminates the influence of a random starting partition and stabilizes the clusterization process.

3. Some Biological Aspects of Cluster Classification

In this section we compare partitions obtained above to certain types of biologically meaningful bacteria classifications. We are going to discuss the biological meaning of the obtained results

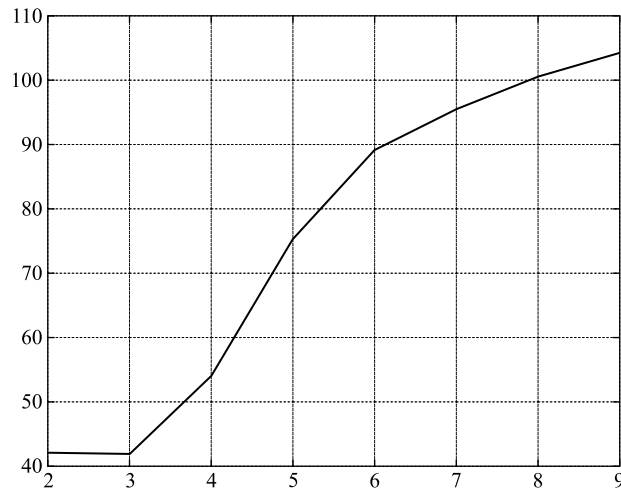


Fig. 2. Graph of the normalized quantities of the misclassified items

from the point of view of their biological relevance. Specifically, we compare partitions built on different alphabets and different genomes parts in order to point out the connection between these clustering solutions and several genomes parameters. We test four types of the distance on set 185 of genomes. Each one of them has been divided into coding and non-coding parts and for each part two spectra have been calculated. The 4-letters alphabet (A, T, C, G) and the reduced 2-letters alphabet ($A = G, T = C$) are used. In Table 1 the quantities of the clusters, in the final optimal partition, are presented.

Table 2 shows that the portion of the misclassified items is very small, approximately 12%. This fact demonstrates high agreement

Table 1

Partition of the set of 185 bacteria based on the standard 4 letters alphabet, for the coding (cod) and the non-coding (non) parts and the partition based on the two-letter (purine/pyrimidine) alphabet for the coding (pcod) and for the non-coding parts

	1	2	3
Cod	53	99	36
Non	47	105	36
Pcod	63	75	50
Pnon	80	48	60

Table 2

Intersections of the partitions, obtained based on coding (cod) and non-coding (non) genome parts, for the four-letter alphabet

	1 non	2 non	3 non
1 cod	45	0	8
2 cod	0	96	3
3 cod	2	9	25

between the obtained partitions and their biological meaning. We think that the significant similarity among partitions, based on the coding and the non-coding genomes parts, using the 4 letters alphabet, is very interesting. In Table 3 the characteristics of the purine/pyrimidine alphabet are exhibited.

Table 3

Intersection of the partitions, obtained based on coding (pcod) and non-coding (pnon) genome parts, for the two-letter alphabet

	1 pnon	2 pnon	3 pnon
1 pcod	50	5	8
2 pcod	30	43	2
3 pcod	0	0	50

In this case, cluster 2, of pcod, includes cluster 2 of pnon and a significant part of the species of cluster 1 of pnon. We observe approximately 24 % of misclassified items. However, the Chi-square values are 233.42 and 195.16 for Tables 2 and 3, accordingly. The suitable P-values are less than 0.001, in both cases, hence, the correlation between the partitions are significant at a level of 0.001. I.e. the presented correlations are not random.

Thus, we concentrate only on the following cases:

- the coding part of the 4-letter alphabet;
- the coding and non-coding part the 2-letter alphabet.

The corresponding partitions are presented in Appendix. We would like to note that the partition's comparison exhibits a considerable dissimilarity between the solutions.

The properties of the obtained clusters are checked with respect to two ecological parameters (oxygen and temperature) and Archaea/Bacteria classification. In order to estimate the compatibility of the parameters with the partitions, we also consider random

partitions. Namely, we start from a particular partition and then we randomly distribute the species, keeping the number of the sets and the number of elements in each set constant. This procedure is repeated 10,000 times. Having obtained the simulated sample we compare the actual factor values to the simulated null hypothesis distribution. Evidence on the non randomness of a factor value is obtained when an actual point is located outside the simulated area.

Oxygen. Bacteria are divided into three groups based on their reaction to oxygen. *Aerobic* bacteria need oxygen for their continued growth and existence. The second group — *Anaerobic* bacteria cannot absorb oxygen. The third group consists of the **Facultative Anaerobes**, which prefer growing in the presence of oxygen, but can continue to grow without it. Each cluster may be characterized by a three-dimensional vector (x, y, z) of the relative frequencies of the bacteria belonging to each of these types. It lays on a simplex $x + y + z = 1$ that makes possible to consider only two dimensional projections of the sets, say (y, z) . Therefore, the test results are presented on the Anaerobic-Facultative Anaerobic plane (Fig. 3).

Figure 3 shows that partitions created by the coding sequences are not random for both alphabets. I.e. appropriate points are located outside of the simulated area. However, for the partitions built on the non-coding sequences the null hypothesis, on the partition randomness, is not rejected. Therefore, it may be assumed that the effect of the oxygen-consumption factor on the DNA sequence is manifested at the level of the coding part of the genome. However, the non-coding genome part appears to be independent of this parameter in the case of the purine-pyrimidine alphabet. This conclusion is supported by the observed random bacteria partition, which is based on this part of the genome.

Temperature. An important characteristic of a bacterium is the temperature at which the bacterium grows most rapidly. The bacteria growing at a fast pace under heat are called thermopiles or hyper-thermopiles. The others are called — mesophilic. In Table 4 we present results of the permutation test on for the temperature parameter.

The frequencies of the thermopiles according to

- (a) the coding parts genomes at the 4-letter alphabet;
- (b) the coding parts genomes at the 2-letter alphabet;
- (c) the non-coding parts genomes at the 2-letter alphabet.

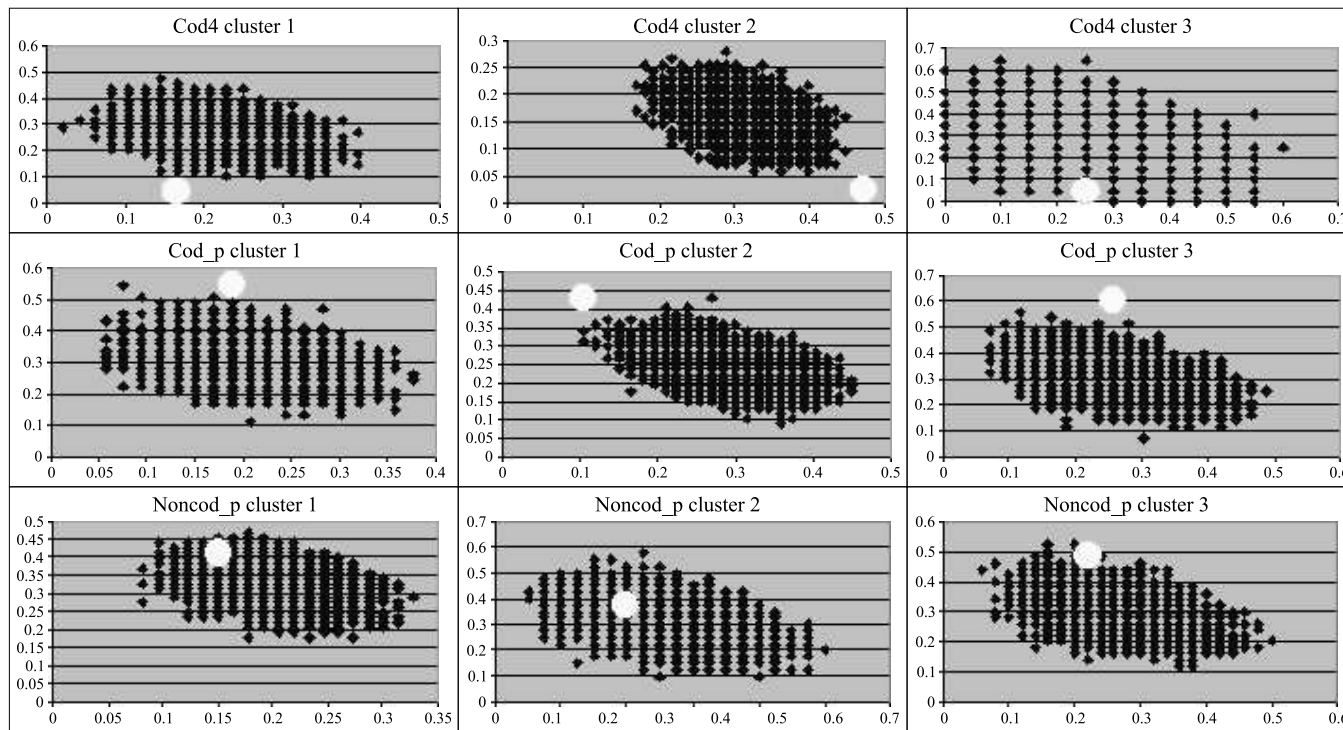


Fig. 3. Results of the permutation test applied for an estimation of the non-randomness of the percent of anaerobic organisms for each of the clusters. The X axis represents the Anaerobic percent and the Y axis represents the Facultative Anaerobic. The allocated circle corresponds to the actual partition

Table 4

Results of the permutation test applied for an estimation of the non-randomness of the percent of the thermophiles for each of the clusters

Cluster	Actual value	min	Max	< a. v.	> a. v.
(a)					
1	0.104	0.021	0.313	3562	6438
2	0.129	0.011	0.172	8827	1173
3	0.133	0.067	0.533	5148	4852
(b)					
1	0	0	0.255	5	9995
2	0.045	0.015	0.224	149	9851
3	0.347	0.02	0.286	10000	0
(c)					
1	0.03	0.015	0.224	40	9960
2	0.024	0	0.317	259	9741
3	0.288	0.017	0.254	10000	0

The number of clusters is indicated in the first column. The second column presents the relative frequencies of the thermophilics in the cluster (actual value (a.v.)). The simulated extreme values are given in the third and fourth columns. The columns named “< a. v.” and “> a. v.” exhibit the quantities of the cases for which the simulated value is less or more than the actual one, correspondently. It can be concluded, that the partition based on the 4 letter alphabet is practically not correlated with the temperature factor. Oppositely, the 2 letter alphabet based partition demonstrates high correlation with the mentioned factor for the coding and the non-coding parts of the genome.

Archaea classification. In this step, a permutation test has been applied to compare the Archaea classification of the Kingdom Bacteria. As it is shown in Table 5, the Archaea set is not detected by means of the 4-letter alphabet. However, the calculations based on the two-letter alphabet exhibit a significant non-randomness of the Archaea distribution. This result corresponds to the ones depicted in our recent publications [15, 18].

Table 5

Results of the permutation test applied for the estimation of the non-randomness of the Archaea percent. The columns' headings in Table 5 are identical to those of Table 4

Cluster	Actual value	min	Max	< a. v.	> a. v.
(a)					
1	0.096	0	0.288	4076	5924
2	0.141	0.03	0.192	9234	766
3	0.045	0	0.455	2567	7433
(b)					
1	0	0	0.254	1	9999
2	0.041	0.014	0.205	123	9877
3	0.347	0	0.265	10000	0
(c)					
1	0	0.025	0.213	10000	0
2	0.065	0	0.283	2196	7804
3	0.288	0	0.237	10000	0

4. Discussion

In this study we present a new method for genome sequences classification using a novel two-step procedure. In the first step, pair-wise comparison of elements is performed on the basis of a certain appropriate functional. Generally speaking, the obtained numerical data are not distances from the point of view of standard metrical requirements. At the second stage each element is matched with all other elements. It has been found that, regardless of the alphabet used, and of the functional significance of the sequence, the obtained number of clusters is the same in all partitions. It has also been shown (see Table 2) that, in the case of the 4-letter alphabet, based on the CS of coding and non-coding parts, the cluster structures are rather similar. Our result suggests that the pair-wise distances for the coding and non-coding parts of the genome for the 4-letter alphabet are more “consistent”, than those for the 2-letter alphabet.

Therefore, the direct use of clustering methods, which are often applicable only for metric spaces, is improper. At the second stage of the algorithm, we match each element of the considered set

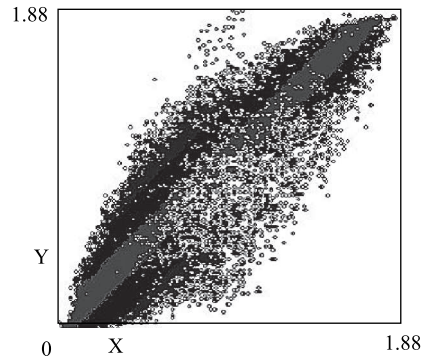


Fig. 4. The correspondence between the distances based on coding and non-coding genome parts in the case of the 4-letter alphabet. Axis X: pair-wise distances based on the coding genome parts. Axis Y — the same, for the non-coding parts

with the vector of its pseudo-distances to all the other elements. This vector will be considered to be an element of a metrical space, which validates the use of the cluster analysis, the desired technique of the pair-wise data estimation being retained. In this work, the pair-wise comparison is based on the CS of sequences and performed by calculating the Spearman coefficient, which is not a metric. The resulting vectors (defined above) are used as elements of the Euclidian space for further clustering. It has been found that, regardless of the alphabet used and of the functional significance of the sequence, the optimal number of clusters is 3 in each partition. It has also been shown (see Table 2) that the cluster structures, based on the CS of coding and non-coding parts of genomes are rather similar to each other in the case of the 4-letter alphabet. In accordance with this, the corresponding distances based on the Spearman coefficient, also correlate with each other (Fig. 4). Our result suggests that the pair-wise distances for the coding and non-coding parts of the genome for the 4-letter alphabet are more “consistent”, than those for the 2-letter alphabet, in the latter case the cluster structures based on the coding and non-coding parts being quite different (Table 3).

The greater difference between the cluster structures for the coding and non-coding genome parts in the case of the 2-letter alphabet in contrast to the virtual coincidence of these structures in the case of the 4-letter alphabet is not significant from the viewpoint of a formal criterion, being, however, substantiated by the biologically significant cluster analysis. Indeed, we have shown that the

pair-wise distances between the coding parts of bacteria correlate with the level of oxygen consumption regardless of the alphabet chosen, while the pair-wise distances between the non-coding parts for the 2-letter alphabet are independent of this parameter. Perhaps, this is connected to the mechanisms the protein resistance to “oxidative stress”. On the other hand, the bacteria’s temperature dependence does not correlate, with genome distances, only for the 4-letter alphabet. The same phenomenon follows from the analysis of the Archae/Bacteria genome classification. The obtained result is well coordinated with the biological context of the problem even though partitions have been received without resorting to biological criteria.

References

1. Rocha E.P., Viari A., Danchin A. Oligonucleotide bias in *Bacillus subtilis*: general trends and taxonomic comparisons // Nucleic Acids Research. 1998. V. 26. P. 2971–2980.
2. Rogozin I.B., Makarova K.S., Natale D.A., Spiridonov A.N., Tatusov R.L., Wolf Y.I., Koonin E. V. Congruent evolution of different classes of non-coding DNA in prokaryotic genomes // Nucleic Acids Research. 2002. V. 29. P. 4264–4271.
3. Holmquist G.P.. Evolution of chromosome bands: molecular ecology of noncoding DNA // J. Mol. Evol. 1989, V. 28. P. 469–486.
4. Baldi P., Basnee P.-F. Sequence analysis by additive scale: DNA structure for sequences and repeats of all lengths // Bioinformatics. 2000. V. 16. P. 865–889.
5. Trifonov E.N. The multiple codes of nucleotide sequences // Bull. Math. Biol. 1989. V. 51. P. 417–432.
6. Li Y. C., Korol A. B., Fahima T., Beiles A., Nevo E. Microsatellites: genomic distribution, putative functions and mutational mechanisms: a review // Mol. Ecol. 2002. V. 11. P. 2453–2465.
7. Snel B., Huynen M. A., Dutilh B. E. Genome trees and the nature of genome evolution // Annu. Rev. Microbiol. 2005. V. 59. P. 191–209.
8. Bussemaker H.J., Li H., Siggia E.D. Building a dictionary for genomes: identification of presumptive regulatory sites by statistical analysis // Proc. Natl. Acad. Sci. U.S.A. 2000. V. 97. P. 10096–10100.
9. Healy J., Thomas E. E., Schwartz J. T., Wigler M. Annotating large genomes with exact word matches // Genome Research. 2003. V. 13. P. 2306–2315.

10. *Robins H., Krasnitz M., Barak H., Levine A.J.* A relative-entropy algorithm for genomic fingerprinting captures host-phage similarities // *J. Bacteriol.* 2005. V. 187. P. 8370–8374.
11. *Qi J., Wang B., Hao B.-I.* Whole Proteome Prokaryote Phylogeny Without Sequence Alignment: A K-String Composition Approach // *Journal of Molecular Evol.* 2004. V. 58. P. 1–11.
12. *Chapus C., Dufraigne C., Edwards S., Giron A., Fertil B., Deschavanne P.* Exploration of phylogenetic data using a global sequence analysis method // *BMC Evolutionary Biology.* 2005. V. 9. P. 63.
13. *Kirzhner V., Korol A., Bolshoy A., Nevo E.* Compositional spectrum — revealing patterns for genomic sequence characterization and comparison // *Physica A.* 2002. V. 312. P. 447–457.
14. *Kirzhner V., Nevo E., Korol A., Bolshoy A.* One promising approach to a large-scale comparison of genomic sequences // *Acta Biotheoretica.* 2003. V. 51. P. 73–89.
15. *Kirzhner V., Paz A., Volkovich Z., Nevo E., Korol A.* Different clustering of genomes across life using the A-T-C-G and degenerate R-Y alphabets: Early and late signaling on genome evolution? // *J. Mol. Evol.* 2007. V. 64. P. 448–456.
16. *Dudoit J., Fridlyand S.* A prediction-based resampling method to estimate the number of clusters in a dataset // *Genome Biology.* 2002. V. 3.
17. *Roth V., Lange V., Braun M., Buhmann J.* Stability-based validation of clustering Solutions // *Neural Computation.* 2004. V. 16. P. 1299–1323.
18. *Kirzhner V., Bolshoy A., Volkovich Z., Korol A., Nevo E.* Large scale genome clustering across life based on a linguistic approach // *BioSystem.* 2005. V. 81. P. 208–222.
19. *Volkovich Z., Kirzhner V., Bolshoy A., Korol A., Nevo E.* The Method of N-grams in Large-Scale Clustering of DNA texts // *Pattern Recognition.* 2005. V. 38. P. 1902–1912.
20. *Kendall M.G.* Rank Correlation Methods. — Charles Griffin & Co, Ltd, London, 1970.
21. *Jain A., Dubes R.* Algorithms for Clustering Data. — Englewood Cliffs, N. Y., Prentice-Hall, 1988.
22. *Gordon A.D.* Classification. — Chapman and Hall, CRC, Boca Raton, FL, 1999.
23. *Chakravarthy S.V., Ghosh J.* Scale-based clustering using the radial basis function Network // *IEEE Transactions on Neural Networks.* 1996. V. 7. P. 1250–1261.

24. Tibshirani R., Walther G., Hastie T. Estimating the number of clusters via the Gap Statistic // J. Royal Statistic. Soc. B. 2001. V. 63. P. 411–423.
25. Rand W. Objective criteria for the evaluation of clustering methods // J. Amer. Statist. Assoc. 1971. V. 66. P. 846–850.
26. Fowlkes E. B., Mallows C. L. A method for comparing two hierarchical clusterings // J. Amer. Statist. Assoc. 1983. V. 78. P. 553–569.
27. Ben-Hur A., Elisseeff A., Guyon I. A stability based method for discovering structure in clustered data // In: Pacific Symposium on Biocomputing. 2002. P. 6–17.
28. Hubert L., Arabie P. Comparing partitions // J. Classification. 1985. V. 2. P. 193–218.
29. Lange T., Roth V., Braun L. M., Buhmann J. M. Stability-based validation of clustering Solutions // Neural Computation 2004. V. 16. P. 1299–1323.
30. Kuhn H. The Hungarian method for the assignment problem // Naval Research Logistics Quarterly. 1955. V. 2. P. 83–97.

Appendix

The table of all of three partitions. In the first column the name of a bacteria, in the second — taxonomic class. In following three columns in figure number of cluster accordingly in partitions (1) a code part genomes on the basis of 4-letter alphabet and (2) the two-letter alphabet and (3) on the basis of non-coding parts genomes is designated at two-letter alphabet. The table is ordered on the first partition, and inside of these clusters — under names of bacteria.

Aeropyrum pernix K1	Thermoprotei	1	3	3
Agrobacterium tumefaciens chromosome I	Alphaproteobacteria	1	2	1
Agrobacterium tumefaciens chromosome II	Alphaproteobacteria	1	2	1
Agrobacterium tumefaciens str. C58 (Dupont) chr I	Alphaproteobacteria	1	2	1
Agrobacterium tumefaciens str. C58 (Dupont) chr II	Alphaproteobacteria	1	2	1
Bifidobacterium longum NCC2705	Actinobacteria	1	2	2
Bordetella bronchiseptica RB50	Betaproteobacteria	1	2	2
Bordetella parapertusis 12822	Betaproteobacteria	1	2	2
Bordetella pertussis Tohama I	Betaproteobacteria	1	2	2

<i>Bradyrhizobium japonicum</i> USDA 110	Alphaproteobacteria	1	2	2
<i>Brucella melitensis</i> 16M chromosome I	Alphaproteobacteria	1	2	1
<i>Brucella melitensis</i> 16M chromosome II	Alphaproteobacteria	1	2	1
<i>Brucella suis</i> 1330 chromosome I	Alphaproteobacteria	1	2	1
<i>Brucella suis</i> 1330 chromosome II	Alphaproteobacteria	1	2	1
<i>Caulobacter crescentus</i> CB15	Alphaproteobacteria	1	2	1
<i>Chlorobium tepidum</i> TLS	Chlorobia	1	2	1
<i>Chromobacterium violaceum</i> ATCC 12472	Betaproteobacteria	1	2	2
<i>Corynebacterium efficiens</i> YS-314	Actinobacteria	1	2	1
<i>Dechloromonas aromatica</i> RCB	Betaproteobacteria	1	2	1
<i>Deinococcus radiodurans</i> R1 chromosom II	Deinococci	1	2	2
<i>Deinococcus radiodurans</i> R1 chromosome I	Deinococci	1	2	1
<i>Desulfovibrio vulgaris</i> subsp. <i>vulgaris</i> str. Hildenborough	Deltaproteobacteria	1	2	2
<i>Geobacter metallire</i> GS-15	Deltaproteobacteria	1	1	1
<i>Geobacter sulfurre</i> PCA	Deltaproteobacteria	1	1	1
<i>Gloeobacter violaceus</i> PCC 7421	Cyanobacteria	1	2	1
<i>Gluconobacter oxydans</i> 621H	Alphaproteobacteria	1	2	1
<i>Haloarcula marismortu</i> ATCC 43049	Euryarchaeota	1	2	2
<i>Halobacterium salinarum</i> sp. NRC-1	Halobacteria	1	2	2
<i>Leifsonia xyli</i> subsp. <i>xyli</i> str. CTCB07	Actinobacteria	1	2	2
<i>Mesorhizobium loti</i> MAFF303099	Alphaproteobacteria	1	2	2
<i>Methanopyrus kandleri</i> AV19	Methanopyri	1	3	3
<i>Methylococcus capsulatus</i> str. Bath	Gammaproteobacteria	1	2	1
<i>Mycobacterium avium</i> subsp. <i>paratuberculo</i> str. k10	Actinobacteria	1	2	2
<i>Mycobacterium bovis</i> subsp. <i>bovis</i> AF2122/97	Actinobacteria	1	2	2
<i>Mycobacterium leprae</i> TN	Actinobacteria	1	2	2
<i>Mycobacterium tuberculosis</i> CDC1551	Actinobacteria	1	2	2
<i>Mycobacterium tuberculosis</i> H37 Rv	Actinobacteria	1	2	2

Natronomonas pharaonis DSM 2160	Euryarchaeota	1	2	2
Nitrobacter winogradsky Nb-255	Alphaproteobacteria	1	2	1
Nocardia farcinic IFM 10152	Actinobacteria	1	2	2
Pelodictyon luteolum DSM 273	Chlorobia	1	1	1
Propionibacterium acnes KPA171202	Actinobacteria	1	2	2
Pseudomonas aeruginosa PAO1	Gammaproteobacteria	1	2	1
Pseudomonas putida KT2440	Gammaproteobacteria	1	2	2
Pseudomonas syringae pv. tomato str. DC3000	Gammaproteobacteria	1	2	2
Ralstonia solanacearum GMI 1000	Betaproteobacteria	1	2	2
Rhodopseudomonas palustris HaA2	Alphaproteobacteria	1	2	2
Sinorhizobium meliloti 1021	Bacilli	1	2	1
Streptomyces avermitilis MA-4680	Actinobacteria	1	2	2
Streptomyces coelicolor A3(2)	Actinobacteria	1	2	2
Synechococcus sp. WH 8102	Chroococcales	1	2	1
Thermobifida fusca YX	Actinobacteria	1	2	1
Thermus thermophilus HB2	Deinococci	1	3	3
Aquifex aeolicus VF5	Aquificae	2	3	3
Archaeoglobus fulgidus DSM 2661	Archaeoglobi	2	3	3
Bacillus anthracis str. Ames	Bacilli	2	1	1
Bacillus anthracis str. Ames 0581 ('Ames Ancestor')	Bacilli	2	1	1
Bacillus cereus ATCC 14579	Bacilli	2	1	1
Bacillus halodurans	Bacilli	2	1	1
Bacillus subtilis subsp. subtilis str. 168	Bacilli	2	1	1
Bacteroides thetaiotaom VPI-5482	Bacteroides	2	1	1
Borrelia burgdorferi B31	Spirochaetes	2	3	3
Buchnera aphidicola str. Bp (Baizong pistacia)	Gammaproteobacteria	2	1	3
Buchnera aphidicola str. Sg (Schizap graminum)	Gammaproteobacteria	2	3	3

<i>Buchnera aphidiocola</i> str. APS (<i>Acyrtosiphon pisum</i>)	Gammaproteobacteria	2	1	3
<i>Campylobacter jejuni</i> subsp. <i>jejuni</i> NCTC 11168	Epsilonproteobacteria	2	3	3
<i>Candidatus Blochmanni floridanus</i>	Gammaproteobacteria	2	1	2
<i>Carboxydotherrmus hydrogenoformans</i> Z-2901	Clostridia	2	3	3
<i>Chlamydia muridarum</i> Nigg	Chlamydiae	2	3	3
<i>Chlamydia trachomatis</i> D/UW-3/CX	Chlamydiae	2	3	3
<i>Chlamydomphila caviae</i> GPIC	Chlamydiae	2	3	3
<i>Chlamydomphila pneumoniae</i> AR39	Chlamydiae	2	3	3
<i>Chlamydomphila pneumoniae</i> CWL029	Chlamydiae	2	3	3
<i>Chlamydomphila pneumoniae</i> J138	Chlamydiae	2	3	3
<i>Clostridium acetobutylicum</i> ATCC 824	Clostridia	2	3	3
<i>Clostridium perfringens</i> st. 13	Clostridia	2	3	3
<i>Colwellia psychrerythraea</i> 34H	Gammaproteobacteria	2	1	1
<i>Desulfotalea psychrophila</i> LSv54	Deltaproteobacteria	2	1	1
<i>Ehrlichia canis</i> str. Jake	Alphaproteobacteria	2	1	2
<i>Ehrlichia ruminanti</i> str. Gardel	Alphaproteobacteria	2	1	2
<i>Enterococcus faecalis</i> V583	Bacilli	2	1	1
<i>Fusobacterium nucleatum</i> subsp. <i>nucleatum</i> ATCC 25586	Fusobacteria	2	3	3
<i>Haemophilus ducreyi</i> 35000HP	Gammaproteobacteria	2	2	1
<i>Haemophilus influenzae</i> Rd KW20	Gammaproteobacteria	2	1	1
<i>Helicobacter hepaticus</i> ATCC 51449	Epsilonproteobacteria	2	3	3
<i>Helicobacter pylori</i> 26695	Epsilonproteobacteria	2	3	3
<i>Helicobacter pylori</i> J99	Epsilonproteobacteria	2	3	3
<i>Lactobacillus johnsonii</i> NCC 533	Bacilli	2	1	1
<i>Lactococcus lactis</i> subsp. <i>lactis</i> I11403	Firmicutes	2	1	1
<i>Leptospira interrogans</i> serovar Copenhagen chromosome II	Spirochaetes	2	3	3
<i>Leptospira interrogans</i> serovar Copenhagen chromosome I	Spirochaetes	2	3	3

<i>Leptospira interrogans</i> serovar lai str. 56601 chromosome I	Spirochaetes	2	3	3
<i>Leptospira interrogans</i> serovar lai str. 56601 chromosome II	Spirochaetes	2	3	3
<i>Listeria innocua</i> Clip11262	Bacilli	2	1	1
<i>Listeria monocytogenes</i> EGD-e	Bacilli	2	1	1
<i>Listeria monocytogenes</i> str. 4b F2365	Bacilli	2	1	1
<i>Mesoplasma florum</i> L1	Mollicutes	2	1	1
<i>Methanobacterium thermoautotrophicus</i> str. Delta H	Methanomicrobia	2	3	3
<i>Methanococcus jannaschii</i> DSM 2661	Methanococci	2	3	3
<i>Methanococcus maripaludis</i> S2	Methanococci	2	3	3
<i>Methanosarcina acetivorans</i> C2A	Methanomicrobia	2	3	3
<i>Methanosarcina barkeri</i> str. fusaro	Methanomicrobia	2	3	3
<i>Methanosarcina mazei</i> Go1	Methanomicrobia	2	3	3
<i>Mycoplasma gallisepticum</i> R	Mollicutes	2	1	1
<i>Mycoplasma genitalium</i> G-37	Mollicutes	2	1	1
<i>Mycoplasma mobile</i> 163K	Mollicutes	2	3	3
<i>Mycoplasma mycoides</i> subsp. <i>mycoides</i> SC str. PG1	Mollicutes	2	1	3
<i>Mycoplasma penetrans</i> HF-2	Mollicutes	2	1	1
<i>Mycoplasma pneumoniae</i> M129	Mollicutes	2	1	1
<i>Mycoplasma pulmonis</i> UABCTIP	Mollicutes	2	3	3
<i>Nanoarchaeum equitans</i> Kin4-M	Nanoarchaeota	2	3	3
<i>Nostoc</i> sp. PCC 7120	Nostocales	2	1	3
<i>Oceanobacillus ihayensis</i> HTE831	Bacilli	2	1	1
Onion yellow phytoplasma OY-M	Mollicutes	2	1	1
<i>Parachlamydia</i> sp. UWE25	Chlamydiae	2	3	3
<i>Pasteurella multocida</i> subsp. <i>multocida</i> str. Pm70	Gammaproteobacteria	2	1	1
<i>Photobacterium profundum</i> SS9 chromosome I	Gammaproteobacteria	2	2	2

Photobacterium profundum SS9 chromosome II	Gammaproteobacteria	2	2	2
Photorhabdus luminescens subsp. laumondii TTO1	Gammaproteobacteria	2	2	1
Picrophilus torridus DSM 9790	Thermoplasmata	2	3	3
Prochlorococcus marinus subsp. marinus str. CCMP1375	Prochlorophytes	2	2	3
Prochlorococcus marinus subsp. pastoris str. CCMP1986	Prochlorophytes	2	3	3
Psychrobacter arcticus 273-4	Gammaproteobacteria	2	2	2
Pyrococcus abyssi GE5	Thermococci	2	3	3
Pyrococcus furiosus DSM 3638	Thermococci	2	3	3
Pyrococcus horikoshii OT3	Thermococci	2	3	3
Rickettsia conorii str. Malish 7	Alphaproteobacteria	2	1	1
Rickettsia felis URRWXCal2	Alphaproteobacteria	2	1	3
Rickettsia prowazekii str. MadridE	Alphaproteobacteria	2	1	2
Rickettsia typhi str. Wilmington	Alphaproteobacteria	2	1	2
Staphylococcus aureus subsp. aureus Mu50	Bacilli	2	1	1
Staphylococcus aureus subsp. aureus N315	Bacilli	2	1	1
Staphylococcus aureus subsp. aureus str MW2	Bacilli	2	1	1
Staphylococcus epidermidis ATCC 12228	Firmicutes	2	1	1
Staphylococcus saprophyticus subsp. saprophyticus ATCC 15305	Bacilli	2	1	1
Streptococcus agalactiae 2603V/R	Bacilli	2	1	1
Streptococcus agalactiae NEM316	Lactobacillales	2	1	1
Streptococcus mutans UA159	Bacilli	2	1	1
Streptococcus pneumoniae R6	Bacilli	2	1	3
Streptococcus pneumoniae TIGR4	Bacilli	2	1	3
Streptococcus pyogenes M1 GAS (SF370)	Bacilli	2	1	1
Streptococcus pyogenes MGAS315	Bacilli	2	1	1
Streptococcus pyogenes MGAS8232	Bacilli	2	1	1
Streptococcus pyogenes SSI-1	Bacilli	2	1	1

<i>Streptococcus thermophilus</i> CNRZ1066	Bacilli	2	1	1
<i>Sulfolobus acidocalda</i> DSM 639	Crenarchaeota	2	3	3
<i>Sulfolobus solfataric</i> P2	Thermoprotei	2	3	3
<i>Sulfolobus tokodaii</i> str. 7	Thermoprotei	2	3	3
<i>Thiomicrospira denitrificans</i> ATCC 33889	Epsilonproteobacteria	2	3	3
<i>Treponema denticola</i> ATCC 35405	Spirochaetes	2	3	3
<i>Wolbachia</i> endosymbiunt of <i>Drosophila melanogaster</i>	Alphaproteobacteria	2	1	1
<i>Wolinella succinogenes</i>	Epsilonproteobacteria	2	3	3
<i>Bdellovibrio bacteriovorus</i> HD100	Deltaproteobacteria	3	3	3
<i>Corynebacterium diphtheriae</i> NCTC 13129	Actinobacteria	3	2	1
<i>Corynebacterium glutamicum</i> ATCC 13032	Actinobacteria	3	2	3
<i>Coxiella burnetii</i> RSA 493	Gammaproteobacteria	3	1	1
<i>Erwinia carotov</i> subsp. <i>atrosep</i> SCRI104	Gammaproteobacteria	3	2	2
<i>Escherichia coli</i> CFT073	Gammaproteobacteria	3	2	2
<i>Escherichia coli</i> K-12	Gammaproteobacteria	3	2	2
<i>Escherichia coli</i> O157-H7	Gammaproteobacteria	3	2	2
<i>Escherichia coli</i> O157-H7 EDL933	Gammaproteobacteria	3	2	2
<i>Geobacillus kaustophilus</i> HTA426	Firmicutes	3	1	1
<i>Idiomarina loihiensis</i> L2TR	Gammaproteobacteria	3	2	1
<i>Lactobacillus plantarum</i> WCFS1	Bacilli	3	2	2
<i>Mannheimia succiniciproducens</i> MBEL55E	Gammaproteobacteria	3	1	1
<i>Neisseria meningitidis</i> MC58	Betaproteobacteria	3	2	1
<i>Neisseria meningitidis</i> Z2491	Betaproteobacteria	3	2	3
<i>Nitrosococcus oceani</i> ATCC 19707	Gammaproteobacteria	3	1	1
<i>Nitrosomonas europaea</i> ATCC 19718	Betaproteobacteria	3	2	1
<i>Nitrospira multiformis</i> ATCC 25196	Betaproteobacteria	3	1	1
<i>Pelobacter carbinolic</i> DSM 2380	Deltaproteobacteria	3	1	1
<i>Pirellula</i> sp. 1 (<i>Rhodophirellula baltica</i> SH 1)	Planctomycetacia	3	2	1
<i>Porphyromonas gingivalis</i> W83	Bacteroides	3	1	1

<i>Prochlorococcus marinus</i> str. MIT 9313	Prochlorophytes	3	1	1
<i>Pyrobaculum aerophilum</i> IM2	Crenarchaeota	3	3	3
<i>Rhodopirellula baltica</i> SH 1	Planctomycetes	3	2	1
<i>Salmonella enterica</i> subsp. <i>enterica</i> serovar Typhi str. CT18	Gammaproteobacteria	3	2	2
<i>Salmonella enterica</i> subsp. <i>enterica</i> serovar Typhi Ty2	Gammaproteobacteria	3	2	2
<i>Salmonella typhimurium</i> LT2	Gammaproteobacteria	3	2	2
<i>Shewanella oneidensis</i> MR-1	Gammaproteobacteria	3	2	1
<i>Shigella flexneri</i> 2a str. 2457T	Alphaproteobacteria	3	2	2
<i>Shigella flexneri</i> 2a str. 301	Alphaproteobacteria	3	2	2
<i>Thermococcus kodakaraensis</i> KOD1	Euryarchaeota	3	3	3
<i>Yersinia pestis</i> biovar <i>Mediaeva</i> str. 91001	Gammaproteobacteria	3	2	2
<i>Zymomonas mobilis</i> subsp. <i>mobilis</i> ZM4	Alphaproteobacteria	3	1	1

EXPERIMENTALLY OBTAINED MEASUREMENTS FOR RECONSTRUCTION OF ORTHOGRAPHIC VIEWS

Miri Weiss-Cohen

Ort Braude College of Engineering, Karmiel, Israel

Alina Bondarenko

Rafael, Haifa, Israel

Yoram Halevi

Technion-Israel Institute of Technology, Haifa, Israel

The orthographic views are constructed from a moving sensor, either a distancemeter or a camera, in an automatic fashion using estimation techniques. The raw distance measurements are processed via a filter which generates estimates of the part dimensions and position. The fact that the outcome of the estimation (measurement) process is a set of explicit contour equations is suitable for creating orthographic views using variational geometry representation.

1. Introduction

The work deals with creating an automatically generated 2D orthogonal views. Theoretically not only that one can reconstruct the full model directly from measurements, but both the object and the sensor may be moving with. However, a typical scenario is a stationary object and measurements taken from orthogonal directions. The measurements can be in-plane laser, orthogonal laser or a camera, and all three cases are discussed.

The main idea in creating the views is that the object consists of primitives with known shape but with unknown parameters, i. e. a cylinder with unknown dimensions. Other quantities that are unknown in the process are the position and orientation of the object. Such problems have been addressed in the past for other purposes, e.g. [3], or using different approach [3, 17]. The methodology of the identification in this paper follows the one in [2]. When the measurements are obtained by a camera, a preliminary step is edge detection. Then there exist algorithms that reconstruct position and orientation using identified objects in consecutive images [7].

With the raw measurements provided by either the laser distancemeter or the camera, the construction of the image in the pre-specified shape becomes a non-linear estimation problem. This is one of the most addressed problems with numerous approaches ranging from standard Least Squares, through Gradient Weighted Least Squares [1], to the more robust M-estimators [9]. The most common approach for estimation from a sequence of measurements is the use of the the Kalman Filter, which is the optimal estimator for linear problems, on a linearized model about the prior estimate. This is the Extended Kalman Filter (EKF) [5, 14]. The Iterative Extended Kalman Filter (IEKF) [10] uses repeated linearizations to increase the accuracy. The recent method of Noise Updated Iterative Extended Kalman Filter (IEKF) [2] goes one step further and uses the identified noise in the linearization. This has a strong effect in cases where the noise effect is strongly non-linear.

A number of approaches have been developed over the past decades to interpret user-supplied orthographic views. These approaches are inputs for reconstruction of a 3D object models. The main reconstruction approach is the wireframe — B-rep bottom up approach [4, 8, 12, 13]. We propose an automatic procedure [15,16] for representing the experimental measurements by elements of variational geometry. In particular, the main novelty in the approach is its use of understanding the nature of 2D engineering drawings. This understanding is translated into an actual algorithm by means of topological relations and dimensional scheme analysis.

2. Creating the 2d Orthographic views

We assume that the measurement is planar. One possibility is an in-plane laser measurement shown schematically in Fig. 1. At each measurement instant the sensor is in a known position and orientation, up to small uncertainty which is modeled as noise. Then the distance to the contour of the object is measured, again with a certain noise. In [2] the most general situation where the situation is 3D and the object is moving as well was considered, but here we assume that the problem is planar and that the object is stationary, yet in an unknown position. Treating the position of the object as unknown is appealing from a practical point of view because it eliminates the need for registration and/or initialization.

The general shape of the object, or parts of the object, is assumed to be known, but not its parameters. In body coordinates

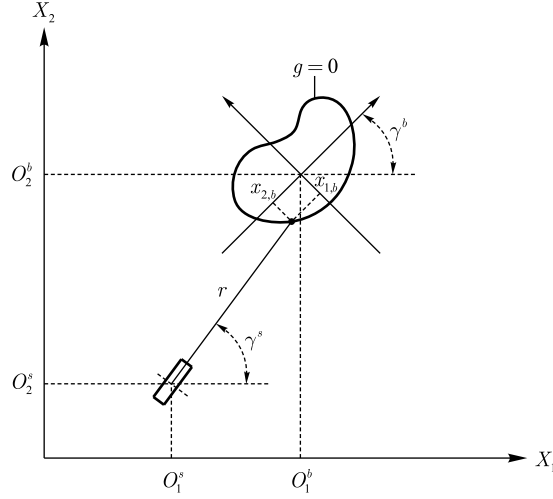


Fig. 1. General setting of the object and the sensor

(see Fig. 1), the contour is given by

$$g(x_1^b, x_2^b, \vartheta) = 0, \quad (2.1)$$

where ϑ is a vector of parameters. For example, in a circle this vector includes the two coordinates of the center and the radius, and in an ellipse it includes the two coordinates of the center and the two radii. g may be a only piecewise continuous or a vector of several functions, e.g. for multi-facet objects, as is the case in Fig. 2. The four functions in that case are

$$\begin{aligned} g_1 &= x_2^b, & g_2 &= x_2^b - B, & g_3 &= x_1^b, & g_4 &= x_1^b - A, \\ \vartheta &= (O_1^b, O_2^b, A, B). \end{aligned} \quad (2.2)$$

Alternatively, the four lines can be described by the single function that is their product

$$g = x_2^b \cdot (x_2^b - B) \cdot x_1^b \cdot (x_1^b - A). \quad (2.3)$$

The advantage of this form is the automatic calculations without the need of logical operations that determines which facet is active. Dealing with simpler functions, on the other hand, is better from a statistical point of view.

To relate the local (2.1) to the actual measurement a coordinate transformation is required. While this can be done by means of standard geometry, the derivation is more structured and pseudo-

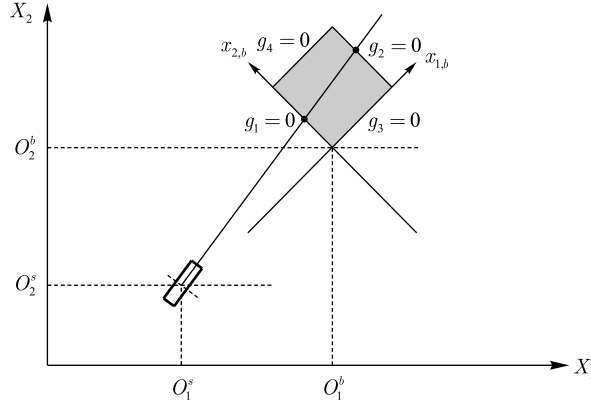


Fig. 2. A rectangle represented by four functions

linear when homogeneous coordinates are used. The body coordinates system and the global coordinates system are related as

$$\begin{bmatrix} X_1 \\ X_2 \\ 1 \end{bmatrix} = T_b R_b \begin{bmatrix} x_{1,b} \\ x_{2,b} \\ 1 \end{bmatrix}, \quad (2.4)$$

where $R_b(\gamma^b)$ and $T_b(O_1^b, O_2^b)$ are the rotation and translation transformation matrices from body coordinates respectively. On the other hand, a point on the contour, from the sensor point of view, is given by

$$\begin{bmatrix} X_1 \\ X_2 \\ 1 \end{bmatrix} = T_s R_s \begin{bmatrix} r \\ 0 \\ 1 \end{bmatrix}, \quad (2.5)$$

where $R_s(\gamma^s)$ and $T_s(O_1^s, O_2^s)$ are the transformations from the sensor coordinates. Combining the two relationships we have

$$\begin{bmatrix} x_{1,b} \\ x_{2,b} \\ 1 \end{bmatrix} = R_b^{-1} T_b^{-1} T_s R_s \begin{bmatrix} r \\ 0 \\ 1 \end{bmatrix}. \quad (2.6)$$

Substitution into the contour (2.1) gives the k -th measurement reading of a point on the contour, in terms of the measured quantities as

$$g(x^s, r, \vartheta) = 0, \quad (2.7)$$

where

$$x^s = [O_1^s \quad O_2^s \quad \gamma^s]^T. \quad (2.8)$$

The derivation so far was purely geometric, assuming perfectly accurate measurements. In reality each of the measured quantities contains noise. The k -th measurement point is described by

$$\underbrace{\begin{bmatrix} x_k^s \\ r_k \end{bmatrix}}_{\text{true}} = \underbrace{\begin{bmatrix} z_k^s \\ \bar{r}_k \end{bmatrix}}_{\text{measured}} - \underbrace{\begin{bmatrix} v_k^s \\ v_k^r \end{bmatrix}}_{\text{noise}} = z_k - v_k, \quad (2.9)$$

where z_k are the actual measurements and the elements of v_k represent the corresponding noises. With this notation (2.7) becomes

$$y_k = g(z_k, \vartheta, v_k) = 0. \quad (2.10)$$

Equation (2.10) is an implicit measurement where the artificial output is always zero but the actual measurements, in particular the distance r , appear as coefficients. The equation is nonlinear in the measurements and consequently in the noise as well. Since there is an uncertainty in the position and the orientation of the sensor, these quantities, which in general change from one measurement to another, need to be estimated in addition to the object parameters. Define the vector

$$\tilde{x}_k = \begin{bmatrix} x_k^s \\ \vartheta \end{bmatrix}. \quad (2.11)$$

The overall measurement is, therefore,

$$\begin{bmatrix} z_k^s \\ 0 \end{bmatrix} = \begin{bmatrix} I & o \\ g(\tilde{x}_k, z_k, v_k) \end{bmatrix} \tilde{x}_k + v^s \quad (2.12)$$

or generically

$$\tilde{y}_k = H(\tilde{x}_k, v_k). \quad (2.13)$$

Notice that although only the parameters ϑ are of interest, they cannot be separated from the rest of \tilde{x}_k one has to estimate x_k^s as well. The implicit measurement is non-linear. One way to overcome that is the use of the extended Kalman filter (EKF), which uses a linearized version of the measurement. For the sake of brevity, we present only a general statement of the estimation scheme, where the details can be found in [3].

$$\hat{\tilde{x}}_k = \tilde{x}_k + K_k(\tilde{y}_k - H(\tilde{x}_k, 0)), \quad (2.14)$$

$$K_k = K_k(\tilde{C}_k, \tilde{D}_k), \quad (2.15)$$

where \tilde{C}_k and \tilde{D}_k are the coefficients of the state and the noise after linearization, which are given by

$$\tilde{C}_k = \begin{bmatrix} I & 0 \\ \frac{\partial g}{\partial x^s} & \frac{\partial g(\tilde{x})}{\partial \vartheta} \end{bmatrix}_{|\tilde{x}=\tilde{x}_k}, \quad \tilde{D}_k = \begin{bmatrix} I & 0 \\ 0 & \frac{\partial g}{\partial r} \end{bmatrix}_{|\tilde{x}=\tilde{x}_k}. \quad (2.16)$$

Since the EKF is based on linearization about the a priori estimation \tilde{x}_k , a natural extension is a recursive procedure with \tilde{x}_k replacing \tilde{x}_k and so on. This is the iterative extended Kalman filter (IEKF) [13]. A further extension is given in [2] with the introduction of noise updated iterative extended Kalman filter (NUIEKF). The key idea is that better estimation of the state variables can be obtained if the measurement noise is updated iteratively as well. In general terms, the estimation in the i -th iteration of the k -th time step is given by

$$\begin{aligned} \hat{x}_{k,i+1} &= \tilde{x}_k + K_{1,k,i}(y_k - H(\hat{x}_{k,i}, \hat{v}_{k,i}) - \tilde{C}_{k,i}(\tilde{x}_k - \hat{x}_{k,i}) + \tilde{D}_{k,i}\hat{v}_{k,i}), \\ \hat{v}_{k,i+1} &= K_{2,k,i}(y_k - H(\hat{x}_{k,i}, \hat{v}_{k,i}) - \tilde{C}_{k,i}(\tilde{x}_k - \hat{x}_{k,i}) + \tilde{D}_{k,i}\hat{v}_{k,i}), \end{aligned} \quad (2.17)$$

where the gains $K_{1,k,i}$ and $K_{2,k,i}$ are calculated based on the iterative linearization $\tilde{C}_{k,i}$, $\tilde{D}_{k,i}$.

Overhead measurements are somewhat simpler as they produce directly points that are, apart from the noise, on the contour $g(x_1^b, x_2^b, \vartheta) = 0$. The estimation process then follows along the same lines as in (2.8)–(2.10) but with (2.10) as the only measurement, i. e. $\tilde{x}_k = \vartheta$ and equations (2.12) and (2.16) are reduced to their lower part.

B. Camera Measurements. Assume now that the measurement is made by a single, or repeated, camera picture. Points belonging to the contour can be found by any edge detection algorithm, so the main task is to find how the coordinates of a point in the 3D global or local space is translated to the 2D image plane. Let (x_c, y_c, z_c) be the coordinates of a point in the camera frame with its Z axis coinciding with the optical axis, and (u, v) the coordinates in the image. The two coordinate systems are related by

$$u = \frac{f_x x_c}{z_c} + u_0, \quad v = \frac{f_y y_c}{z_c} + v_0, \quad (2.18)$$

where f_x, f_y are the focal distances of the camera, and u_0, v_0 are the coordinates of the image center. The same relationship can be also written as

$$\begin{bmatrix} u' \\ v' \\ w \\ 1 \end{bmatrix} = K \cdot \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & u_0 & 0 \\ 0 & f_y & v_0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}, \quad (2.19)$$

where the (u', v') and (u, v) are related by

$$u = \frac{u'}{w}, \quad v = \frac{v'}{w}. \quad (2.20)$$

A crucial point in camera measurement is the missing depth information, i.e. while u and v are given in the picture, w is not. Hence additional information is required to transformation back to (x_c, y_c, z_c) . Let $G(\tilde{x})$ describe the surface of the 3D object, as shown in Fig. 3. The line connecting a point \tilde{x}_1 , which will be on the contour in the picture, and the camera position O_c is tangent to the object, or equivalently, perpendicular to the normal. Mathematically this is described by

$$\frac{\partial G}{\partial \tilde{x}} \Big|_{\tilde{x}=\tilde{x}_1} \perp (\tilde{x}_1 - O_c) \Rightarrow \left(\frac{\partial G}{\partial \tilde{x}} \Big|_{\tilde{x}=\tilde{x}_1} \right)^T \cdot (\tilde{x}_1 - O_c) = 0. \quad (2.21)$$

If the body has sharp edges, e.g. a box, the situation becomes more complex, and the single normal is replaced by a cone that the line should be perpendicular to a certain member of it. Since we focus on a orthogonal views, a more plausible assumption is that the distance from the object described by w , is known. With that information the scaling from the image to the true view is straightforward. As was mentioned earlier, edge detection methods are used to define points that, apart from noise, lie on the contour. The situation then is similar to the overhead laser measurements, and the estimation problem is fitting a known function with unknown parameters.

At the end of the process, same as the case of laser measurement, we have two or three orthogonal views. In the next section we outline the procedure for creating a 3D object out of them.

3. Topological relations and dimensioning analysis using variational geometry

The input to this stage is a dimensioned 2D view, which goes through a constraint evaluation process resulting in a 2D view constraint set. Dimensions define geometric constraints, such as distance between two points, distance between a point and a bar, and an angle between two bars. Spatial relations define topological constraints such as tangency, parallelism, and perpendicularity [11]. The constraints extracted from each 2D view represent relations among explicit and implicit characteristics. Each dimension is formulated as a constraint. There are two kinds of constraints, one, defined by a single equation, and a compound constraint, which require two or more equations.

Each constraint equation is a function of points in the geometric dimension scheme. Equation i , denoted f_i , is formulated as follows:

$$f_i = \{x_1, y_1, x_2, y_2, \dots, x_n, y_n\}, \quad (3.1)$$

where n denotes the number of points constraints by the geometric entity.

For the complete 2D view, a set of constraints, denoted as F , is given as follows:

$$F = \{f_1, f_2, \dots, f_m\}. \quad (3.2)$$

As an example, a distance from a point to a line is presented. To constrain the distance D between a point P_a and line P_bP_c two vectors must be defined: a unit vector \hat{U} from P_b to P_c and a vector \bar{V} from P_b to P_a the distance D is a cross product $\hat{U} \times \bar{V}$:

$$\hat{U} = \frac{x_c - x_b}{|P_bP_c|}\hat{i} + \frac{y_c - y_b}{|P_bP_c|}\hat{j} = U_x\hat{i} + U_y\hat{j} \quad (3.3)$$

and

$$\bar{V} = (x_b - x_a)\hat{i} + (y_b - y_a)\hat{j}, \quad (3.4)$$

where i and j are unit vectors in the x and y directions, respectively.

The point-to-line constraint is formalized as:

$$f_1 = U_x(y_b - y_a) - U_y(x_b - x_a) - D = 0. \quad (3.5)$$

The system constructs a knowledge base of variational geometry rules for constraining the dimensioning scheme and the relations between the geometry sites in the view [16]. In the rule base, constraining the dimensioning scheme is done by positioning a selected point, called anchor point, at the origin in order to

prevent solid body translation. All points are defined relative to this anchor point. To prevent solid body rotation, a bar is defined to be horizontal, i. e., parallel to the x axis. For the dimensioning and the constraint set to be valid, the Jacobian constraint matrix should meet two requirements. First, the number of constraints must be equal to twice the number of vertices, and second, the rank of the matrix must equal the number of constraints. Meeting these requirements indicates that the matrix is non-singular and hence there is neither redundancy nor lack of dimensions and definitions

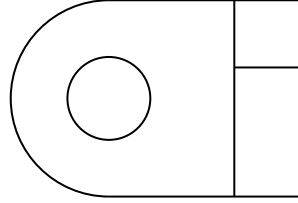


Fig. 3. Front view

of the constraints. The constraint set $F\{\text{Front}\}$, for the Front view in Fig. 3 is formalized in (3.6):

$$\begin{aligned}
 f_1: & (x_2 - x_1)^2 + (z_2 - z_1)^2 - c^2 = 0 \quad \text{Euclidian Distance,} \\
 f_2: & (x_3 - x_3)^2 + (z_3 - z_3)^2 - d^2 = 0 \quad \text{Euclidian Distance,} \\
 f_3: & (x_1 - x_8)^2 + (z_1 - z_8)^2 - c^2 = 0 \quad \text{Euclidian Distance,} \\
 f_4: & (x_8 - x_6)^2 + (z_8 - z_6)^2 - b^2 = 0 \quad \text{Euclidian Distance,} \\
 f_5: & (x_9 - x_{10})^2 + (z_9 - z_{10})^2 - R_2^2 = 0 \quad \text{Euclidian Distance,} \\
 f_6: & (x_9 - x_4)^2 + (z_9 - z_4)^2 - R_1^2 = 0 \quad \text{Euclidian Distance,} \\
 f_7: & (x_1 - x_8)(x_2 - x_1) + (z_1 - z_8)(z_2 - z_1) = 0 \quad \text{Perpendicularity,} \\
 f_8: & (x_2 - x_3)(x_9 - x_3) + (z_2 - z_3)(z_9 - z_3) = 0 \quad \text{Perpendicularity,} \\
 f_9: & (x_1 - x_2)(x_7 - x_2) + (z_1 - z_2)(z_7 - z_2) = 0 \quad \text{Perpendicularity,} \\
 f_{10}: & (x_5 - x_6)(x_8 - x_6) + (z_5 - z_6)(z_8 - z_6) = 0 \quad \text{Perpendicularity,} \\
 f_{11}: & (x_5 - x_4)(x_9 - x_4) + (z_5 - z_4)(z_9 - z_4) = 0 \quad \text{Perpendicularity,} \\
 f_{14}: & (x_1 - x_8)(x_7 - x_8) + (z_1 - z_8)(z_7 - z_8) = 0 \quad \text{Perpendicularity,} \\
 f_{12}: & x_5 - x_2 = 0 \quad \text{Collinear Points,} \\
 f_{13}: & x_6 - x_1 = 0 \quad \text{Collinear Points,} \\
 f_{15}: & x_4 - x_3 = 0 \quad \text{Collinear Points,}
 \end{aligned}$$

$$\begin{aligned}
f_{16}: \quad & x_{10} - x_3 = 0 \quad \text{Collinear Points,} \\
f_{17}: \quad & x_9 - x_3 = 0 \quad \text{Collinear Points,} \\
f_{18}: \quad & x_1 = 0 \quad \text{Anchor Point,} \\
f_{19}: \quad & z_1 = 0 \quad \text{Anchor Point,} \\
f_{20}: \quad & z_2 - z_1 = 0 \quad \text{Orientation.}
\end{aligned} \tag{3.6}$$

A triangular prism is described in Fig. 4 by two orthographic view: Front view (F) and Side view (S).

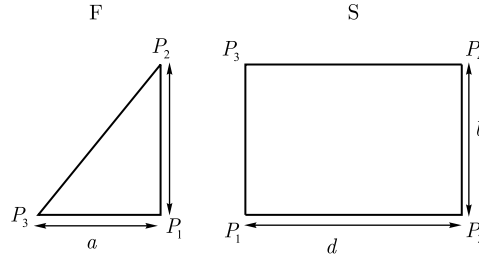


Fig. 4. Front view (F) and Side view (S) of a triangular prism

In this example we use the following parameters for the vertices found in the two 2D views.

For the front view we define the three vertices $P_1(x_1, y_1)$, $P_2(x_2, y_2)$, $P_3(x_3, y_3)$ and for the Side view define the four vertices $P_1(x_1, z_1)$, $P_2(x_2, z_2)$, $P_3(x_3, z_3)$, $P_4(x_4, z_4)$. Two rules from the variational geometry rule base described in Section 3 are used:

Rule #1 — Euclidian distance between two point, and

Rule #2 — Perpendicularity.

To prevent solid body translation, the anchor point was fixed as the origin (0,0). All the points are defined relative to this anchor point. To prevent solid body rotation, we choose a particular bar to be horizontal, i.e., parallel to the horizontal axis.

For the Front view of Fig. 5 we formulate the following six constraints:

$$\begin{aligned}
f_1: \quad & (x_2 - x_1)^2 + (y_2 - y_1)^2 - b^2 = 0, \\
f_2: \quad & (x_3 - x_1)^2 + (y_3 - y_1)^2 - a^2 = 0, \\
f_3: \quad & (x_1 - x_2)(x_3 - x_1) + (y_1 - y_2)(y_3 - y_1) = 0, \\
f_4: \quad & x_1 = d = 0,
\end{aligned}$$

$$\begin{aligned}
f_5: \quad & x_1 = e = 0, \\
f_6: \quad & y_2 - y_1 = 0.
\end{aligned} \tag{3.7}$$

From these set of equations we obtain a 6x6 Jacobian matrix. The rank of the matrix was calculated and the dimensioning was found to be proper and equal twice the number of vertices.

For the Side view, eight other constraints are formulated as follows:

$$\begin{aligned}
f_1: \quad & (y_2 - y_1)^2 + (z_2 - z_1)^2 - d^2 = 0, \\
f_2: \quad & (y_3 - y_1)^2 + (z_3 - z_1)^2 - b^2 = 0, \\
f_3: \quad & (y_4 - y_3)^2 + (z_4 - z_3)^2 - d^2 = 0, \\
f_4: \quad & (y_4 - y_2)^2 + (z_4 - z_2)^2 - b^2 = 0, \\
f_5: \quad & (y_1 - y_2)(y_3 - y_1) + (z_1 - z_2)(z_3 - z_1) = 0, \\
f_6: \quad & y_1 = e = 0, \\
f_7: \quad & z_1 = a = 0, \\
f_8: \quad & z_2 - z_1 = 0.
\end{aligned} \tag{3.8}$$

As before, we calculated an 8×8 Jacobian matrix from the set of equations. The rank of this matrix was found to be 8, indicating proper dimensioning and constraint definition of the side view as well.

4. Conclusions

A comprehensive method for automatically constructing a parametric representation for an orthographic views obtained by a moving sensor was reported and described. It is assumed that the shape of the object body, actually sub objects of the entire body, is known but its parameters are not. Also, the position of the body need not be known in advance. The raw distance measurements are processed via a filter which generates, estimates of the part dimensions and position. Since the formulation leads to implicit measurement equations, standard extended Kalman filter techniques, usually fail to converge to accurate values. A new method, called a Noise Updated Iterative Extended Kalman Filter, was developed and used.

The fact that the outcome of the estimation (measurement) process is a set of explicit contour equations is suitable to the second step which is creating the orthographic views through variational geometry representation. In previous applications those equations had to be built from the views as a preliminary step. The match between the output of the Kalman filtering approach for 2D reconstruction, and the starting point for the parametric representation and orthographic views reconstruction is the key advantage for the integrated approach. The output representation of the orthographic views are the input for the process of reconstruction of 3D models [6].

References

1. Beck J. V., and Arnold K. J. Parameter Estimation in Engineering and Science. — Wiley series in probability and mathematical statistics. J. Wiley, N. Y., 1977.
2. Bondarenko A., Halevi Y., and Shpitalni M. Object identification and tracking via noise updated iterative extended Kalman filter. — 7th Biennial ASME Conference on Engineering Systems Design and Analysis (ESDA), Manchester, UK, 7 p., 2004.
3. Borenstein J., and Koren Y., Obstacle avoidance with ultrasonic sensors // IEEE Journal of Robotics and Automation. 1988. V. RA-4, No. 2. P. 213–218.
4. Chen K.-Z., and Feng X.-A., Holo-extraction of information from paper drawing for 3D reconstruction // Computer Aided Design. 2002. V. 34. P. 665–677.
5. De Geeter J., Van Brussel H., and De Schutter J. A smoothly constrained Kalman filter // IEEE Trans. on Pattern Analysis and Machine Intelligence. 1997. V. 19, No. 10.
6. Dori D., and Weiss M. A scheme for 3D object reconstruction from dimensioned orthographic views // Engineering Applications of Artificial Intelligence. 1996. V. 9. P. 52–64.
7. Faugeras O. D. Three Dimensional Vision, A Geometric Viewpoint. — MIT press, 1993.
8. Geng W., and Wang J. Embedding visual cognition in 3D reconstruction from multi-view engineering drawing // Computer Aided Design. 2002. V. 34. P. 321–336.
9. Huber P. J. Robust Statistics. — John Wiley & Sons, N. Y., 1981.
10. Jazwinsky A. M. Stochastic Processes and Filtering Theory. — Academic, N. Y., 1970.

11. *Lin V.C., Light R.A., and Gossard D.C.* Variational geometry in computer aided design // *Computer Graphics*. 1981. V. 14. P. 171–177.
12. *Liu S., Hu S., Chen Y., and Sun J.* Reconstruction of curved solids from engineering drawing // *CAD*. 2001. V. 33. P. 1059–1072.
13. *Markovsky G., and Wesley M.A.* Fleshing out wireframes // *IBM J. of Research and Development*. 1980. V. 24, No. 5. P. 582–597.
14. *Mendel J.M.* *Lessons in Digital Estimation Theory*. — Prentice-Hall, 1987.
15. *Weiss M., and Dori D.* Variational geometry as a tool for dimension validation of recognized 2D views in engineering drawings, shape, structure and pattern recognition // *World Scientific*. 1995. P. 416–431.
16. *Weiss-Cohen M.* Reconstruction of Solids Models from Orthographic views using knowledge retrieval and composite graphs // *Computer-Aided Design and Applications (CAD&A)*. 2007. V. 4. P. 159–167.
17. *Zhang Z.* Parameter estimation techniques: A tutorial with application to conic fitting // *Image and Vision Computing Journal*. 1997. V. 15, No. 1. P. 59–76.

ON INTRUSION IN A LINEARLY-STRATIFIED AMBIENT: THE ASYMMETRIC STEADY-STATE MODEL

Tamar Zemach

Department of Computer Science, Technion, Haifa

Marius Ungarish

Software Engineering Department, ORT Braude College
of Engineering, Karmiel, Israel

The behavior of the steady intrusive gravity currents spreading into a stratified ambient fluid is investigated. The intrusive gravity current of thickness h and density ρ_c which propagates with speed U at the neutral buoyancy level of a long horizontal channel of height H into a stratified ambient fluid whose density increases linearly from ρ_o to ρ_b is investigated. The intrusive and the ambient fluids are assumed to be asymmetric due to axis passing the stagnation point of the system. The Boussinesq, high-Reynolds number two-dimensional configuration is discussed. The Long's model combined with the flow-force balance over the width of the channel and the pressure balances over a density current are used to obtain the desired results. It is shown that the intrusion velocity decreases with decreasing the asymmetry of the system and approaches its minimum for the symmetric configuration. In addition, the comparison between asymmetric and symmetric configurations shows no significant differences between the models.

1. Introduction

A gravity current is formed whenever one fluid flows primarily horizontally into a lighter or heavier fluid. The classical theoretical work on the subject, Benjamin [3], considers gravity currents entering homogeneous fluid using steady-state theory, though including a model for energy dissipation.

In the recent work, Ungarish [15] generalized the classical results of Benjamin concerning the propagation of a steady boundary gravity current of density into a homogeneous ambient to the case of a stratified ambient. For a Boussinesq, high-Reynolds two-dimensional configuration, a flow-field solution of Long's model, combined with flow-force balance over the width of the channel was

used. In particular, the study shows that for a weak stratification the classical result of Benjamin is fully recovered.

Previous investigation of the asymmetric intrusion was done by Holyer and Huppert [7]. They extended Benjamin's work to currents of prescribed volume flux and constant density ρ_c entering an ambient fluid as either a boundary current or an intrusion. The ambient non-continuously stratified fluid was composed from the two layers: the upper layer of the constant density ρ_{up} and the lower layer of the constant density ρ_{lo} , where $\rho_{up} < \rho_c < \rho_{lo}$. It was shown that the depth of the current is not always uniquely determined and it is necessary to use some additional, to the conservation relationships, assumptions to determine which solution occurs. An appropriate principle is obtained by considering dissipative currents. It was postulated that the energy which is lost will go to form a stationary wave train behind the current.

Additional investigation of the asymmetric intrusion was done by Cheong, Kuenen and Linden [4]. They made a number of experiments and presented numerical results for the propagation of the constant density fluid into the two-layer non-continuously stratified ambient fluid configuration. In particular, it was shown that if the density of the intrusion is the depth-weighted mean of the layer densities, the intrusion propagation speed approaches its minimum value.

The present work attempts to extend the steady-state theory of the boundary gravity currents, presented by Ungarish, to the intrusive gravity currents propagating into an linearly continuously stratified ambient fluid. The system of equations obtained by Ungarish for the boundary gravity currents is quite complicated and in some cases non unique solution is obtained. The theory become much more complicated in the intrusion case.

The structure of the paper is as follows. In Section 2 the problem is formulated in form of the system of two non-linear equations. In Section 3 results are presented and various values of governing parameters are examined and discussed. In Section 4 the asymmetric and symmetric cases are compared. In Section 5 some concluding remarks are given.

2. Formulation

The system configuration is sketched in Fig. 1. The long horizontal channel of height $H = H_1 + H_2$ filled with the linearly

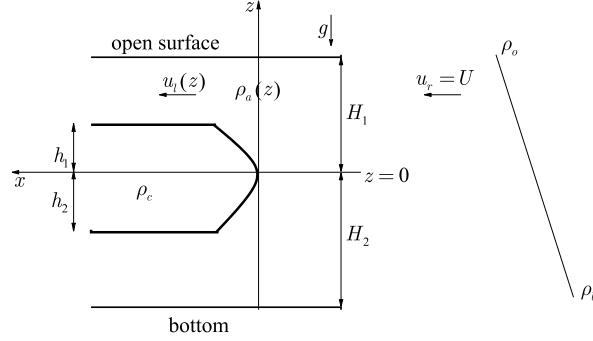


Fig. 1. Schematic description of the system

stratified ambient fluid. The density of the unperturbed ambient, on the right-hand side, increases linearly from ρ_o at the top to ρ_b at the bottom. The denser fluid of density ρ_c and thickness $h = h_1 + h_2$, called the intrusion, propagates with uniform velocity U at the neutral buoyancy level of a stratified fluid. The steady-state flow pattern is concerned.

The driving force is the reduced gravity:

$$g' = \varepsilon g, \quad (2.1)$$

where g is the gravitational acceleration and

$$\varepsilon = \frac{\rho_c - \rho_o}{\rho_o}. \quad (2.2)$$

The horizontal and vertical coordinates are x, z and the corresponding velocity components are u, w . The gravity acts in the $-z$ direction. We choose the $z = 0$ axis passing the stagnation point of the system. Hereafter we will denote by “1” the intrusion part in the $0 \leq z \leq h_1$ domain and the correspond ambient part in the $0 \leq z \leq H_1$ domain. And by “2” another part of intrusion in the $-h_2 \leq z \leq 0$ domain and of the ambient in the $-H_2 \leq z \leq 0$ domain.

An important quantity that characterizes continuously stratified fluids is the buoyancy frequency \mathcal{N} , which is defined by

$$\mathcal{N}^2 = -\frac{g}{\rho_o} \frac{d\rho}{dz} = \frac{g}{H_1 + H_2} \frac{\rho_b - \rho_o}{\rho_o}. \quad (2.3)$$

The other important parameters of the problem are the height relations:

$$a_1 = \frac{h_1}{H_1}; \quad a_2 = \frac{h_2}{H_2}; \quad \eta = \frac{H_1}{H_2}. \quad (2.4)$$

A dimensional parameter which denotes the relative magnitude of the stratification is defined by:

$$S = \frac{\rho_b - \rho_o}{\rho_c - \rho_o}. \quad (2.5)$$

The domain of interest is $S > 1$. We note that $0 \leq S \leq 1$ case corresponds to the boundary gravity current discussed by Ungarish [15] and will not be discussed here.

We denote by Fr the Froude number which is defined by

$$\hat{U} = Fr = \frac{U}{\sqrt{0.5 g'(h_1 + h_2)}}. \quad (2.6)$$

The symmetric case is obtained, if $h_1 = h_2$ and $H_1 = H_2$. In this case the Froude number is defined by $Fr = \frac{U}{(g'h_1)^{1/2}}$ and, as it will be shown later, $S = 2$.

The analysis of the system is started with a solution of a two-layer stratified flow-field over a rigid topography in a channel with an upper and lower horizontal rigid lids at $z = H_1$ and $z = H_2$ correspondingly. The obstacle (or topography) encountered by the unperturbed stratified fluid is defined by the elevation function $z = \chi_1(x)$ in the upper layer and by $z = \chi_2(x)$ in the lower layer. In the $x < 0$ domain, $\chi_1(x) = \chi_2(x) = 0$ and for $x > 0$, far downstream at the left, a parallel geometry is achieved with $\chi_1(x) = h_1 = \text{const} > 0$ and $\chi_2(x) = -h_2 = \text{const} < 0$.

2.1. Upstream flow. The far upstream flow at the right $x \rightarrow -\infty$ consists of parallel horizontal streamlines with constant velocity U and a prescribed stable linearly changing density. Using the subscript “ r ” (right) to denote this region and employing hydrostatic balance, we obtain:

$$\begin{aligned} u_r(z) &= U, & -H_2 \leq z \leq H_1, \\ \rho_r(z) &= \frac{\rho_o - \rho_b}{H_1 + H_2} z + \frac{\rho_o H_2 + \rho_b H_1}{H_1 + H_2}, & -H_2 \leq z \leq H_1, \\ p_r(z) &= -g \int_0^z \rho_r(z) dz = -g \left[\frac{\rho_o - \rho_b}{H_1 + H_2} \frac{z^2}{2} + \frac{\rho_o H_2 + \rho_b H_1}{H_1 + H_2} z \right], & -H_2 \leq z \leq H_1, \end{aligned} \quad (2.7)$$

where u is the velocity; ρ is the density and p is the pressure.

2.2. Downstream flow. Under the assumption of a two-dimensional, steady Boussinesq hydrostatic flow, Long's model (see Baines [2]) can be applied to each layer separately to reduce the set of the governing Navier-Stokes equations to a single ODE equation for each layer for the displacement $\delta(x, z)$ with proper boundary conditions. In particular, the boundary conditions are: 1). $\delta = 0$ at the upstream right region; 2). $\delta = \chi_1(x)$ at the upper layer and 3). $\delta = -\chi_2(x)$ at the lower layer.

For the ambient above the upper intrusion layer, $h_1 \leq z \leq H_1$, the displacement $\delta(x, z)$ satisfies the following problem

$$\begin{cases} \delta_{zz} + k^2 \delta = 0, \\ \lim_{x \rightarrow -\infty} \delta(x, z) = 0, \\ \delta(x, z = H_1) = 0, \\ \delta(x, z = h_1(x)) = h_1(x), \end{cases} \quad (2.8)$$

where $k = \frac{N}{U}$. The exact solution of (2.8) is given by

$$\delta(z) = h_1(x) \cdot \frac{\sin\left(kH_1 \left[1 - \frac{z}{H_1}\right]\right)}{\sin\left(kH_1 \left[1 - \frac{h_1}{H_1}\right]\right)}. \quad (2.9)$$

By the similar way, for the ambient below the lower intrusion layer $-H_2 \leq z \leq -h_2$,

$$\delta(z) = -h_2(x) \cdot \frac{\sin\left(kH_2 \left[1 + \frac{z}{H_2}\right]\right)}{\sin\left(kH_2 \left[1 - \frac{h_2}{H_2}\right]\right)}. \quad (2.10)$$

We replace the solid obstacle with a stationary fluid of density ρ_c in the domain $x \leq 0$, $-\chi_2(x) \leq z \leq \chi_1(x)$. The xz system is now a frame of reference attached to the gravity current and the origin is the front stagnation point. By combining these two layers we have:

$$\delta(z) = \begin{cases} \frac{h_1}{\sin \gamma_1} \cdot \sin \left[\frac{\gamma_1}{1 - a_1} \left(1 - \frac{z}{H_1}\right) \right], & h_1 \leq z \leq H_1, \\ -\frac{h_2}{\sin \gamma_2} \cdot \sin \left[\frac{\gamma_2}{1 - a_2} \left(1 + \frac{z}{H_2}\right) \right], & -H_2 \leq z \leq -h_2, \end{cases} \quad (2.11)$$

where

$$\frac{\gamma_1}{H_1(1-a_1)} = \frac{\gamma_2}{H_2(1-a_2)} = k = \frac{\mathcal{N}}{U} \quad (2.12)$$

and

$$\gamma_1 = (1-a_1) \sqrt{S \frac{\eta}{\eta+1}} \cdot \sqrt{\frac{2\eta}{a_1\eta+a_2}} \frac{1}{\widehat{U}}, \quad (2.13)$$

In the left region the parallel flow satisfies:

$$\rho_l(x, z) = \begin{cases} \rho_c, & -h_2 \leq z \leq h_1, \\ \rho_r(z) - \frac{\rho_o - \rho_b}{H_1 + H_2} \cdot \frac{h_1}{\sin \gamma_1} \cdot \sin \left[\frac{\gamma_1}{1-a_1} \left(1 - \frac{z}{H_1} \right) \right], & h_1 \leq z \leq H_1, \\ \rho_r(z) + \frac{\rho_o - \rho_b}{H_1 + H_2} \cdot \frac{h_2}{\sin \gamma_2} \cdot \sin \left[\frac{\gamma_2}{1-a_2} \left(1 + \frac{z}{H_2} \right) \right], & -H_2 \leq z \leq -h_2, \end{cases} \quad (2.14)$$

and

$$u_l(x, z) = \begin{cases} 0, & -h_2 \leq z \leq h_1, \\ U \left(1 + \frac{a_1}{1-a_1} \cdot \frac{\gamma_1}{\sin \gamma_1} \cdot \cos \left[\frac{\gamma_1}{1-a_1} \left(1 - \frac{z}{H_1} \right) \right] \right), & h_1 \leq z \leq H_1, \\ U \left(1 + \frac{a_2}{1-a_2} \cdot \frac{\gamma_2}{\sin \gamma_2} \cdot \cos \left[\frac{\gamma_2}{1-a_2} \left(1 + \frac{z}{H_2} \right) \right] \right), & -H_2 \leq z \leq -h_2, \end{cases} \quad (2.15)$$

The hydrostatic balance yields

$$p_l(x) = p_S - g \int_0^z \rho_l(z) dz, \quad (2.16)$$

where S denotes the stagnation point. By Bernoulli's low at the stagnation point, we have

$$p_S = \frac{1}{2} U^2 \rho_r(z=0) = \frac{1}{2} U^2 \frac{\rho_o H_2 + \rho_b H_1}{H_1 + H_2}.$$

2.3. The flow-force balance. Following Benjamin [3] and Ungarish [15], we consider the momentum balance in a fixed rectangular control volume whose lower and upper boundaries are the

zero-stress planes $z = H_1$ and $z = -H_2$ and the vertical boundaries are in the parallel up- and down-stream regions. The assumption of steady-state impose the flow-force balance over the width of the channel

$$\int_{-H_2}^{H_1} (\rho_l u_l^2 + p_l) dz = \int_{-H_2}^{H_1} (\rho_r u_r^2 + p_r) dz. \quad (2.17)$$

The evaluation of the integral of the momentum flux is simplified by the Boussinesq assumption $\rho_{l,r} u_{l,r}^2(z) \approx \rho_o u_{l,r}^2(z)$. After some algebra and more use of the Boussinesq assumption, we can express the flow-force balance as

$$\frac{\widehat{U_1^2}}{2} (M_f(a_1, \gamma_1) + \eta \cdot M_f(a_2, \gamma_2)) = F_1(\gamma_1, a_1) + F_2(\gamma_2, a_2), \quad (2.18)$$

where

$$M_f(a, \gamma) = \frac{1}{1-a} [1 + a - 2a^2 + a^2(\gamma^2 + (\gamma \operatorname{ctg} \gamma)^2 + \gamma \operatorname{ctg} \gamma)]; \quad (2.19)$$

$$F_1(\gamma_1, a_1) = 1 - 0.5 a_1 + \frac{H_1}{H_1 + H_2} \times S \left(-1 + a_1 - \frac{1}{3} a_1^2 + (1 - a_1)^2 \frac{1 - \gamma_1 \operatorname{ctg} \gamma_1}{\gamma_1^2} \right) \quad (2.20)$$

and

$$F_2(\gamma_2, a_2) = \frac{a_2}{a_1 \eta^2} \cdot \left[-1 + 0.5 a_2 + \frac{H_2}{H_1 + H_2} \times S \left(0.5 a_2 + \frac{H_1}{H_2} (1 - 0.5 a_2) - \frac{1}{3} a_2^2 + (1 - a_2)^2 \frac{1 - \gamma_2 \operatorname{ctg} \gamma_2}{\gamma_2^2} \right) \right]. \quad (2.21)$$

The right-hand side of (2.18) can be regarded as the buoyancy pressure driving, and the left-hand side as the dynamic reaction. It is noted that there is no mixing between the upper- and lower-layers terms in both hands of (2.18).

We note, that substitution of $a_2 = a_1$ and $H_1 = H_2$ into (2.18) brings it to the flow-force balance for the symmetric intrusion. In this case, as expected, the equation (2.18) becomes identical to this obtained by Ungarish [15] for the boundary gravity current.

2.4. Pressure balance. The additional assumption should be done in form of the pressure balance between the upper- and lower-layers of the intrusion at the left side of the system.

According Bernoulli's low and using the z -hydrostatic pressure distribution we obtain the condition, which connects between the layers:

$$\eta a_1 \operatorname{ctg} \gamma_1 = \pm a_2 \operatorname{ctg} \gamma_2 \quad (2.22)$$

Finally, by substitution (2.13) into (2.18) and (2.22), we obtain the system of two non-linear equations with four parameters: S , γ_1 , a_1 and a_2 . However, only one pair of these parameters are independent. In Section 3 we will find the solution of this system, (γ_1, a_2) , as function of S and a_1 .

2.5. Validity-stability and criticality. The system of the equations (2.18), (2.22) may have non-unique solution (γ_1, a_2) . The values of γ_1 (or γ_2) may be larger than $\pi/2$. This introduces the possibility of negative u_l . Let us denote by

$$\vartheta_i = \begin{cases} 0 & \left(0 < \gamma \leq \frac{\pi}{2}\right), \\ \frac{a_i}{1-a_i} \gamma_i |\operatorname{ctg} \gamma_i| & \left(\frac{\pi}{2} < \gamma < \pi\right), \\ \frac{a_i}{1-a_i} \frac{\gamma_i}{|\sin \gamma_i|} & (\pi < \gamma), \end{cases} \quad (2.23)$$

where $i = 1, 2$, a measure of the most severe relative negative contribution of the perturbation flow to the resulting u_l . According Baines [2] and Ungarish [15], the results are physically acceptable only for $(\vartheta_1, \vartheta_2) < (1, 1)$ case. The valid results are in the range $0 \leq \gamma_i \leq \max[\pi, (1-a_i)/a_i]$, $i = 1, 2$.

In the following analysis we are concerned only with "valid" solutions, for which $(\vartheta_1, \vartheta_2) < (1, 1)$ holds. Moreover, in this case the \pm sign in equation (2.22) vanishes and it yields

$$\eta a_1 \operatorname{ctg} \gamma_1 = a_2 \operatorname{ctg} \gamma_2. \quad (2.24)$$

3. Results

We assume that the density of the ambient fluid at the stagnation point is equal to the density of the current, $\rho_a(z=0) = \rho_c$, and therefore

$$S = 1 + \frac{1}{\eta}. \quad (3.1)$$

Our first numerical experiment of the model is to keep constant the height of the upper layer of the ambient, H_1 , and the height of the intrusion upper layer, h_1 (or the upper layer height relation $a_1 = h_1/H_1$) and to increase the height H_2 of the ambient of the lower layer. In this case the stratification does not change and so the value of \mathcal{N} . Valid results were obtained in the range $1.0 \leq H_2/H_1 \leq 1.27$. Typical valid results of Fr and the height relation of the lower layer a_2 , obtained numerically from (2.18) and (2.24) are presented in Fig. 2 (see supplementary sheet 2). The case $H_2/H_1 = 1$ corresponds to the symmetric configuration $H_1 = H_2$. As expected, for this case the value of a_2 is equal to a_1 and the values of Fr are equal to the values of Fr mentioned by Ungarish [15] for the boundary gravity current.

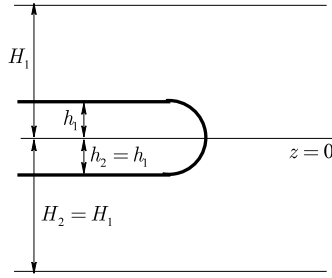
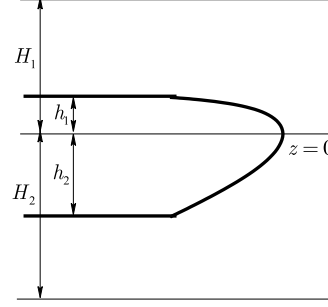
The results show that if the height of the lower layer, H_2 , increases, the lower height ratio a_2 also increases. Moreover, in this case the lower height ratio a_2 becomes greater than a_1 independently on the values of a_1 .

The schematic sketch of this experiment is presented in Fig. 3 for $a_1 = 0.25$: figure (a) shows the symmetric system configuration with $H_1 = H_2 = 1$ and $a_1 = a_2 = 0.25$; figure (b) shows the correspond asymmetric configuration: the upper layers remains unchanged with $a_1 = 0.25$ and the height H_2 of the lower layer increases to be $H_2 = 1.25 H_1$. As result of it, the height h_2 increases and the height ratio a_2 becomes equal to ≈ 0.54 .

The second numerical experiment considers the opposite configuration: it keeps the values of H_2 and a_2 constant and increases the height of the upper layer H_1 . The results show that the upper layer behaves exactly like the lower layer in the first numerical experiment described above. This result is expected although there is no symmetry between the variables a_1 and a_2 in equation (2.18). The valid results were obtained in the range $1.0 \leq H_1/H_2 \leq 1.27$.

From the first two numerical experiments following that the valid results are obtained only in the range $0.786 \leq H_2/H_1 \leq 1.27$. Otherwise, the height relation a_i ($i = 1, 2$) approaches not physical values which are greater than 0.5.

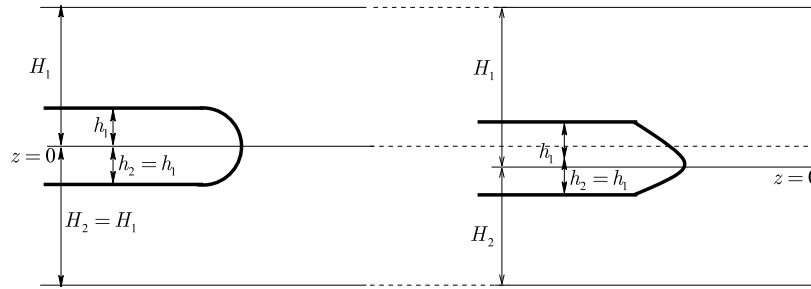
As an additional, third, experiment of the model the value of the ambient height ratio H_2/H_1 is prescribed constant and the value of dimensionless Froude number Fr is calculated as function of the

(a) Symmetric case: $H_2/H_1 = 1$; $a_1 = 0.25$ (b) Asymmetric case: $H_2/H_1 = 1.25$; $a_1 = 0.25$ **Fig. 3.** The first numerical experiment: H_1 and h_1 are fixed and $a_1 = 0.25$; The lower layer height H_2 is increased to $H_2/H_1 = 1.25$

total height ratio a defined by

$$a = \frac{h_1 + h_2}{H_1 + H_2} = \frac{ka_1 + a_2}{k + 1}. \quad (3.2)$$

The sketch of this experiment is presented in Fig. 4: (a) shows the symmetric configuration with $H_2/H_1 = 1$ and $a = 0.25$; in (b) the total height H and a remain constant (with $a = 0.25$), but the height ratio H_2/H_1 increases to be 1.25.

(a) Symmetric case: $H_2/H_1 = 1$; $a = 0.25$ (b) Asymmetric case: $H_1/H_2 = 1.25$; $a = 0.25$ **Fig. 4.** Third numerical experiment: a is fixed. The symmetric $H_2/H_1 = 1.0$ and $a = 0.25$ and its corresponding asymmetric configuration with $H_1/H_2 = 1.25$ and $a = 0.25$

The results are shown in Fig. 5 (see supplementary sheet 2). The solid line corresponds to the symmetric configuration: $a_1 = a_2$ and $H_2/H_1 = 1$. The other lines are the graphs of the

$H_1/H_2 = 0.95, 0.9, 0.85, 0.8$. The Froude number Fr is a decreasing function of a . In the symmetric case, $H_1 = H_2$, Fr approaches its minimal values. Then, when the height ratio H_1/H_2 increasing and the system configuration becomes asymmetric, the Froude number increasing also. This effect was also denoted by Cheong, Kuenen and Linden [4] for the two-layer non-continuously stratified ambient fluid. In particular, it was shown numerically and experimentally that if the density of the intrusion is the depth-weighted mean of the layer densities, the intrusion propagation speed approaches its minimum value.

Following our numerical results, the system (2.18), (2.22) has unique validity solution (γ_1, a_2) for $0 < a_1 \leq 0.5$ and $0.786 \leq H_2/H_1 \leq 1.27$.

However, according Ungarish [15], for the boundary gravity currents, non-unique solutions may be obtained for $0 < a \leq 0.1$. The only case in which this claim is wrong is $S_U \rightarrow 1$, where only unique solution is reported.

As it was discussed above, $S_U \rightarrow 1$ case corresponds to the $S \rightarrow 2$ intrusion case for which unique solution is obtained, which is in agreement with the boundary gravity current results. When $2 < S < 2.3$ and $a_1 \leq 0.1$, the value of a_2 increases and approaches values which are greater than 0.1, so the validity conditions does not satisfied for the lower layer. The same situation takes place for $1.75 < S < 2.0$, when the validity condition does not satisfied for the upper layer.

The additional verification of the model was done by the pressure difference analysis which show that the pressure differences, as expected, approaches small values.

4. Comparison with symmetric case

To sharpen the insights provided by our results obtained for the asymmetric case, we compare to (specially created) corresponding symmetric configurations.

Hereafter all parameters of the symmetric configuration will be denoted by the upper index “(s)”. The total height of the ambient for the symmetric case is set to be equal to the total height $H = H_1 + H_2$ of the ambient for the asymmetric case H and $H_1^{(s)} = H_2^{(s)} = \frac{H}{2}$. The height of the symmetric intrusion

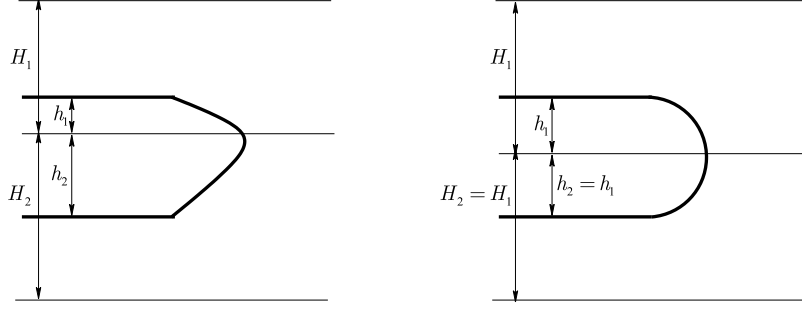
(a) Asymmetric case: $H_2/H_1 = 1.25$; $a_1 = 0.25$ (b) Symmetric case: $H_2/H_1 = 1.0$; $a_1 = 0.39$ 

Fig. 6. Comparison between asymmetric and symmetric configurations: H and h are fixed. The asymmetric $H_2/H_1 = 1.25$ and $a_1 = 0.25$ and its correspond symmetric configuration with $H_2/H_1 = 1$ and $a_1 = 0.39$

is equal to the total intrusion height $h = h_1 + h_2$ of the asymmetric case: $h_1^{(s)} = h_2^{(s)} = \frac{h}{2}$. The typical comparison is sketched on Fig. 6: (a) displays the typical configuration for $H_2/H_1 = 1.25$ and $a_1 = 0.25$ and (b) is a corresponding symmetric configuration with $H_2/H_1 = 1.0$ and $a_1 = a_2 = 0.39$.

Figure 7 (see supplementary sheet 2) shows the behavior of the non-dimensional velocity of the intrusion as a function of the ambient height relation H_2/H_1 . In this case the heights of the asymmetric upper ambient, H_1 , and the intrusion, h_1 , layers are unchanged. The lower layer height increases. The velocity U was scaled by constant value $\mathcal{N}H_1$ and was plotted for various values of a_1 .

The symmetric intrusion height relation was calculated by
$$a^{(s)} = \frac{h}{H} = \frac{\eta a_1 + a_2}{\eta + 1}.$$

We see that the differences between the symmetric and the asymmetric configurations are not large. For $a_1 \leq 0.2$ case, the symmetric intrusion propagates a little faster ($\approx 1.5\%$) than the asymmetric intrusion, however this tendency is changed and for the $a_1 \geq 0.3$, the asymmetric intrusion is faster than the symmetric one.

Similar conclusion was obtained for the two-layer (non-stratified) model of Cheong etc. [4]: the velocity of the intrusion has its minimum, if the two ambient layers have the same height

(symmetric case). This case corresponds to $0.3 \leq a \leq 0.5$, since this is a possible physical range of the stratified intrusion problem, after release.

5. Summary

The steady-state theory of the boundary gravity currents propagating into an linearly-stratified ambient fluid have been generalized to the intrusive gravity currents propagate at the neutral buoyancy level of a long horizontal channel into a stratified ambient fluid. The two-dimensional asymmetric configuration due to the neutral buoyancy level was investigated. The problem was formulated using the flow-force and pressure balances of the system and solved numerically.

The study shows the following points: 1). the valid results was obtained in the range $0.786 \leq H_2/H_1 \leq 1.27$. 2). the velocity of intrusion propagation approaches its minimum in the symmetric configuration case. 3). there is no significant differences between asymmetric and corresponding symmetric configurations.

However, the lack of experimental data prevents sharper conclusions about the insights provided by this theory. We hope that the present study will provide the background and motivation for the laboratory experiments on this problem.

References

1. *Amen R. and Maxworthy T.* The gravitational collapse of a mixed region into a linearly stratified fluid // *J. Fluid Mech.* 1984. V. 96. P. 65–80.
2. *Baines P.* Topographic effects in stratified flows. — Cambridge Univ. Press, 1995.
3. *Benjamin T.* Gravity currents and related phenomena // *J. Fluid Mech.* 1968. V. 31. P. 209–248.
4. *Cheong H., Kuenen J. and Linden P.* The front speed of intrusive gravity currents // *J. Fluid Mech.* 2006. V. 552. P. 1–11.
5. *Faust K. and Plate E.* Experimental investigation of intrusive gravity currents entering stably stratified fluids // *J. Hydraulic Res.* 1984. V. 22. P. 315–325.
6. *Hoult D.* Oil spreading on the sea // *Annu. Rev. Fluid Mech.* 1972. V. 2. P. 341–368.

7. *Hoyler J. and Huppert H.* Gravity currents entering a two-layer fluid // *J. Fluid Mech.* 1980. V. 100. P. 739–767.
8. *Huppert H. and Simpson J.* The slumping of gravity currents // *J. Fluid Mech.* 1980. V. 99. P. 785–799.
9. *Maxworthy T., Leilich J., Simpson J. and Meiburg E.* The propagation of a gravity currents in a linearly stratified fluid // *J. Fluid Mech.* 2002. V. 453. P. 371–394.
10. *Morton K. W., Mayers D. F.* Numerical solution of PDE. — Cambridge Univ. Press, 1998.
11. *de Rooij F.* Sedimenting particle-laden flows in confined geometries. — PhD thesis, DAMTP, University of Cambridge, 1999.
12. *Rottman J. and Simpson J.* Gravity currents produced by instantaneous release of a heavy fluid in a rectangular channel // *J. Fluid Mech.* 1983. V. 135. P. 95–110.
13. *Simpson J.* Gravity currents in the environment and the laboratory. — Cambridge Univ. Press, 1997.
14. *Stommel H. and Farmer H. G.* Abrupt change in width in two-layer open channel flow // *J. Mar. Res.* 1952. V. 11. P. 205–214.
15. *Ungarish M.* On gravity currents in a linearly stratified ambient: a generalization of Benjamin's steady-state propagation results // *J. Fluid Mech.* 2006. V. 548. P. 49–68.
16. *Ungarish M.* Intrusive gravity currents in a stratified ambient-shallow-water theory and numerical results // *J. Fluid Mech.* 2005. V. 535. P. 287–323.
17. *Ungarish M. and Huppert H.* On gravity currents propagating at the base of a stratified ambient // *J. Fluid Mech.* 2002. V. 458. P. 283–301.
18. *Ungarish M. and Huppert H.* On gravity currents propagating at the base of a stratified ambient: effects of axial symmetry and rotation // *J. Fluid Mech.* 2004. V. 521. P. 69–104.
19. *Ungarish M. and Zemach T.* On axisymmetric intrusive gravity currents in a stratified ambient — shallow-water theory and numerical results // *European Journal of Mechanics — B/Fluids.* March–April 2007. V. 26, No. 2. P. 220–235.
20. *Zatsepin A. and Shapiro W.* A study of axisymmetric intrusions in a stratified fluid // *Izvestia Akademii Nauk SSSR, Fizika Atmosfery i Okeana.* 1982. V. 18, No. 1. P. 77–80.
21. *Zemach T.* Gravity currents: two-layer and asymptotic extensions. — M. Sc. Thesis, Technion, Israel, 2002.

ABSTRACTS (in Russian)

Об имитационном подходе в задаче кластерной стабильности

З. Барзилай, М. Голани, З. Волькович

В данной статье обсуждается новый подход к задаче определения числа кластеров в заданной совокупности. Наш метод сочетает в себе методику, основанную на оценке плотности совокупности, с методикой кластерной стабильности. Следуя первой методике, мы рассматриваем кластеры как «острова высокой плотности» в «море данных низкой плотности». Кроме того мы полагаем, что эти острова устойчивы по отношению к зашумлению данных. Иначе говоря, мы считаем, что добавление подходящего шума к данным не ведет к резкому изменению кластеров. С целью проверки кластерной устойчивости мы рассматриваем пары выборок таким образом, что одна из выборок выбирается из рассматриваемой совокупности, а вторая получается с помощью добавления случайного шума к первой. Расстояния между выборками измеряются на основе простых вероятностных метрик, являющихся, де факто, статистиками тестов однородности. Наиболее сконцентрированное в нулевой точке эмпирическое распределение такой статистики соответствует правильному выбору числа кластеров. Численные эксперименты демонстрируют высокую надежность предлагаемого метода.

Классификация моделей локализации неисправностей компьютерных систем

С. Френкель, Е. Левнер, В. Захаров

В статье предлагается классификация алгоритмов поиска, используемых в задачах локализации неисправностей в сложных вычислительных системах. Полезность такой классификации состоит в том, что она позволяет при разработке процедур обслуживания соответствующей системы выбирать алгоритм поиска, оптимальный (в смысле минимизации некоторой функции стоимости) для данной системы и характеристик ее обслуживания. Данная классификация строится аналогично классификации Кенделла в теории очередей. Основой для классификации

является предложенная в предыдущих работах авторов характеристика концептуальной модели систем тестирования и локализации неисправностей, и ряд теорем об оптимальности алгоритмов поиска.

Сложность и состоятельность статистических критериев

А. А. Грушо, Н. А. Грушо, Е. Е. Тимонина

В статье рассматривается связь состоятельности статистических критериев и асимптотической сложности их вычисления в случае конечных пространств. Показано, что из заданной состоятельной последовательности критериев можно построить другую состоятельную последовательность, для которой сложность вычисления принадлежности наблюдаемых значений к критическим множествам этих критериев асимптотически мала по сравнению с аналогичной сложностью для исходных критериев. Однако такое упрощение по сути оказывается фиктивным. Для того, чтобы не допускать фиктивного упрощения вычисления принадлежности наблюдаемых значений к критическим множествам в последовательности критериев, необходимо накладывать дополнительные ограничения на классы рассматриваемых критериев. В статье показано, что в случае естественных ограничений упрощение вычисления может привести к нарушению свойства состоятельности последовательности критериев.

Приведены оценки эффективности использования двухступенчатых критериев, когда сначала работают простые, но не состоятельные в заданном классе альтернатив критерии, и только в случае непринятия гипотезы применяются сложно вычисляемые критерии для всего класса альтернатив.

Байесовские модели обслуживания и надежности

А. А. Кудрявцев, С. Я. Шоргин, В. С. Шоргин, В. М. Ченцов

Рассматривается байесовский подход для определенных задач теории массового обслуживания и теории надежности. Соответствующий метод предусматривает рандомизацию характеристик систем относительно некоторых априорных распределений. Данный подход может использоваться, в частности, для вычисления средних значений и построения доверительных интервалов для вероятностно-временных и надежности характеристик больших групп систем или устройств. Представлены результаты для некоторых моделей параметров входных потоков и времен обслуживания.

Г-сеть с переменной маршрута*Р. Мандзо, А. В. Печинкин*

Рассматриваются сети массового обслуживания с отрицательными заявками (Г-сети), пуассоновским входящим потоком положительных заявок, неэкспоненциальными узлами и зависимым обслуживанием в различных узлах. Каждая заявка, поступающая в сеть, определяется следующими случайными параметрами: длиной маршрута, маршрутом, объемами и временами обслуживания на последовательно проходимых этапах маршрута. Отрицательная заявка при поступлении в сеть может «убить» одну из положительных заявок. Однако «убитая» заявка не покидает сеть, а продолжает путешествовать по сети в соответствии с новым (случайным) маршрутом. Для таких Г-сетей показано, что многомерное стационарное распределение вероятностей состояний сети представимо в мультипликативной форме.

Стохастическая томография с веерной схемой сканирования*О. В. Шестаков*

В томографическом эксперименте, основанном на веерной схеме сканирования, объект облучается расходящимся пучком лучей, испускаемых источником, движущимся вокруг объекта. При использовании такой схемы проекции регистрируются значительно быстрее, чем при использовании традиционной параллельной схемы. В некоторых биологических и физических приложениях исследуемый объект описывается случайной функцией. В работе рассматривается задача восстановления вероятностных характеристик объекта по его проекциям.

Оценка распределения задержки для динамики ВИЧ-инфекции*А. Н. Ушакова*

В настоящей статье рассматриваются два метода оценки распределения задержки в биологических динамических системах. Примером такой системы служит модель ВИЧ-инфекции. Первый метод основывается на параметрическом подходе и на аппроксимации плотности распределения задержки гамма-плотностью. Второй метод является непараметрическим и основан на решении уравнения свертки с выбором параметра регуляризации через параметрический старт.

Об лингвистической классификации геномов бактерий*З. Волькович, В. Кирзнер, З. Барзилай*

Данная статья посвящена классификации 185 полных геномов прокариот на основе модифицированного метода составного спектра. Эта модификация предполагает раздельное вычисление спектра для кодирующей и не кодирующей частей генома. Такой подход позволяет выбрать количество кластеров для классификации геномов, не используя дополнительной информации. Биологический смысл классификации, найденной на основе двух- и четырехбуквенных алфавитов, подтверждает правильность полученных результатов.

Экспериментальные измерения при реконструкции ортографических проекций*М. Вейсс-Козн, А. Бондаренко, Й. Галеви*

Рассматриваются ортографические проекции, получаемые в автоматическом режиме на основе оценочного подхода. Данные измерений фильтруются с помощью оценки измерений размерности и положения. Результатом является ряд явных контурных уравнений, дающих возможность создать, используя вариационные геометрические принципы, ортографические проекции объекта.

Об интрузии в линейно-стратифицированной среде: асимметричная стационарная модель*Т. Цемах, М. Унгарии*

В настоящей статье исследуются проникающие (интрузивные) гравитационные течения, распространяющиеся в стратифицированной жидкости. Рассматриваемые течения, имеющие толщину h , плотность ρ_c и скорость U на уровне нейтральной плавучести, распространяются в канале высотой H , плотность окружающей жидкости в котором растет линейно от ρ_o до ρ_b . Проникающая и окружающие жидкости предполагаются асимметричными относительно вертикальной оси, проходящей через точку стагнации. Рассматривается двумерная система в аппроксимации Буссинеска при больших числах Рейнольдса. Показано, что скорость интрузии (проникновения) убывает при уменьшении асимметрии в системе и достигает минимума в симметричной конфигурации. Сравнение асимметричной и симметричной конфигураций демонстрирует отсутствие принципиальных различий между указанными моделями.

CONTENTS

Preface	3
<i>Barzily Z., Golani M., Volkovich Z.</i> On a simulation aproach to cluster stabilty validation	6
<i>Frenkel S., Levner E., Zakharov V.</i> An approach to classification of computer systems faults localization models.	16
<i>Grusho A., Grusho N., Timonina E.</i> Complexity and Consistency of Statistical Criteria	32
<i>Kudryavtsev A., Shorgin S., Shorgin V., Chentsov V.</i> Bayesian queueing and reliability models	40
<i>Manzo R., Pechinkin A.</i> G-network with the route change	54
<i>Shestakov O.</i> Fan-beam stochastic tomography.	62
<i>Ushakova A.</i> Estimation of delay distribution in HIV dynamics. . . .	78
<i>Volkovich Z., Kirzhner V., Barzily Z.</i> On linguistic classification of bacterial genomes.	86
<i>Weiss-Cohen M., Bondarenko A., Halevi Y.</i> Experimentally Obtained Measurements for Reconstruction of Orthographic Views	112
<i>Zemach T., Ungarish M.</i> On intrusion in a linearly-stratified ambient: the asymmetric steady-state model	125
Abstracts.	139