

of distance learning [8]. A system of algometric algebra has been used in a grammatical analysis for the categorization of text documents and in determining the author's style [9–11], etc.

Of particular importance under conditions of globalization, is the task on identifying the authorship of texts. An analysis of the subject area that we conducted reveals that in most cases the differentiation of phonostatic structures of styles in the process of establishing the author of the text, as well checking a text for plagiarism, involved methods, models, and software tools at the lexical level of a language. However, the phonological level differs from other levels of the language by a stricter structure and ordering of the elements. It is easier to formalize and mathematize. Therefore, it is advisable to apply methods, models, and software tools at the phonological level of a language in order to identify the author of the text and to check the text for plagiarism. Accordingly, the development of methods, models, and tools that would enable the IT differentiation of phonostatic structures of functional styles in the English language is an important and relevant task.

2. Literature review and problem statement

The task on identifying the authorship of the text implies differentiation of texts. Texts are differentiated at the different levels of a language in order to identify their differences and similarities. Thus, the differentiation of texts at the lexical level was performed when modeling grammatical structures [12]. However, the lexical level is an open system. The number of elements is not constant. The system is updated with new words (neologisms) while rarely used words become archaic. An author's style reflects changeable processes in a lexical system. Therefore, the identification of authorship at the lexical level is of a probabilistic character. It is worth noting that grammatical structures are abstract, idealized models, and do not provide for a complete reflection of the speech process. This makes it difficult to define the differential attributes of the author's style. Modeling of semantic structures was used for text differentiation [13]. Semantic structures are the abstract constructs whose implementation depends on the context. That is why a focus on semantics predetermines a probabilistic character of the author's attribution. Texts are differentiated at the lexical and semantic levels when splitting a sentence into key words [14]. Determining the dominant lexical units was used when distinguishing texts in the areas of culture and tourism [15, 16]. Determining the dominant key words does not make it possible to cover lexical vocabulary characteristic of a particular author and is not promising in identifying the author's style. It should be noted that the results of text differentiation at the lexical and semantic levels of a language have a more probabilistic character than that at the phonological level. In contrast to a phonological level, the number of elements is not constant and that compromises the accuracy of calculations. In addition, no combination of the most effective quantitative methods was determined to differentiate texts at each level of a language [17]. When establishing the differential attributes of the author's style using statistical methods, no scheme style→substyle→author was applied, which facilitates determining statistical parameters for the author's manner of presentation in texts from different subjects [18]. Information technologies were not employed

for the author's attribution at the phonological level, and that does not provide the proper level of accuracy [19, 20]. Software systems do not implement a combination of statistical methods, which would provide efficiency of the author's attribution [21]. An analysis of the scientific literature that we conducted revealed that the task on improving the accuracy of text differentiation remains unsolved. To solve the problem, it is required to carry out author's attribution at the phonological level, to apply the combination of statistical methods that is the most efficient to obtain probable results and to determine the degree of validity of factors related to style, substyle, and the author's manner of presentation.

3. The aim and objectives of the study

The aim of present study is to improve the accuracy of differentiation of phonostatic structures of styles in the English language based on the developed methods, models, and software tools for the implementation of the author's, substyle, and style text attribution.

To accomplish the aim, the following tasks have been set:

- to develop a mathematical basis for the system of differentiation of phonostatic structures of functional styles in the English language using the theory of mathematical statistics, which would make it possible to improve the accuracy of output results;
- to construct models for the differentiation of phonostatic structures of styles of the English language;
- to devise a structure of the system and the software that would be based on a modular principle, which would make it possible to rapidly modify the developed IT tools and to ensure that the software system is platform-independent.

4. Development of the system's mathematical basis

The core of any software system is a mathematical basis that includes the developed methods. The constructed mathematical basis for the differentiation of phonostatic structures of styles in the English language includes the following.

1. A method of comprehensive analysis for the differentiation of phonostatic structures of styles [22, 23] is based on the proposed combination of such statistical methods as: a method of hypotheses, a method of ranking, and a method for determining the distances between styles. The algorithm of the constructed method of hypotheses includes the following steps:

Step 1. Check the conformity of frequency of consonant phonemes to the law of normal distribution using the Pearson criterion and a simplified criterion by Romanovsky.

Step 2. Differentiation of texts for the Student's criterion.

Step 3. Determine the groups of consonant phonemes, based on which we established substantial differences in the pairwise comparison of texts.

An algorithm of the ranking method includes the following steps:

Step 1. Determine the mean frequency of groups of consonant phonemes.

Step 2. Construction of descending series of mean frequencies for each group of phonemes.

Step 3. Determine significant differences between the pairwise compared texts based on the difference in ranking.

An algorithm of the method for determining the distances between styles is implemented by the following steps:

Step 1. Differentiation of pairwise compared texts based on the Student's criterion.

Step 2. Derive from the formula for the Student's criterion a formula for determining the distances between styles ($l = \frac{t-t_0}{t}$) [24, 25].

Step 3. Determine a large, medium, and insignificant distance between styles.

The method considered makes it possible to differentiate with greater accuracy the styles, substyles, and texts by different authors.

2. A multi-factor method for determining the degrees of action of the factors related to style, substyle, and the author's manner of presentation, is based on the developed scheme style→substyle→author in order to identify the authorship of texts of the same style, one substyle, but by different authors. An algorithm of the method includes the following basic steps:

Step 1. Determine substantial differences in the pairwise comparison of texts based on the Student's criterion: different styles, different substyles, different authors.

Step 2. Determine a significant, medium, and insignificant degree of action factors related to: style, substyle, the author's manner of presentation.

The method makes it possible to establish with a higher accuracy the affiliation of the text under study to a specific style, substyle, and to identify its author.

5. Development of models for the differentiation of phonostatic structures of styles

Based on the developed methods, we have built statistical models for the style, substyle, and author's differentiation of texts by the ranking method. An algorithm of the specified models includes the following steps.

Step 1. Determine the mean frequency of groups of consonant phonemes for texts: of different styles, different substyles, by different authors, determine the highest and lowest indicators of values for the mean frequency, determine large, medium, and minor differences based on the proposed formula

$$r_{x_1 - x_2}^a = r_{\max x_1}^a - r_{\min x_2}^a.$$

The models developed make it possible to take into consideration, with a greater accuracy, the position of a phoneme in a word, to perform the style, substyle, and the author's attribution of texts based on the ranking difference.

We have developed a statistical model for determining a general stylistic markedness of the examined text. An algorithm for constructing the model includes the following steps:

Step 1. Determine essential differences, based on the Student's criterion, in the compared texts: different styles, different substyles, by different authors, in various subjects.

Step 2. It is proposed to determine the mean value for the three obtained t -values for the Student's criterion:

$$sm = \frac{t_{f_1} + t_{f_2} + t_{f_3}}{3}.$$

Step 3. Determine a large, medium, and insignificant stylistic markedness of the examined text.

The developed model is a combination of three models represented in papers [26, 27]. The model needs to be applied in the case when texts belong to the same style and substyle, but they are by different authors and address a different topic. The model makes it possible to identify the author of texts on various subjects with a higher accuracy. Therefore, the developed methods, models, and algorithms make it possible to improve the accuracy of differentiation of the phonostatic structures of styles.

6. Development of the structure and software of the system

The methods and models developed have been implemented in the programming language java, in the system of differentiation of phonostatic structures of styles in the English language.

The structure of the developed software is shown in Fig. 1; it is based on a modular principle and allows individual customization and support for each module, it ensures high reliability of the system [28]; the built software is easily upgraded.

The algorithm the English language style differentiation based on the mean frequencies of groups of consonant phonemes, which is implemented in the system, implies the execution of a sequence of the following basic steps:

1. Computer processing of the examined text:
 1. 1. Upload an English-language text to the software.
 1. 2. Convert the text into a transcription variant.
 1. 3. Separate from the transcription characters those that denote consonant phonemes.
 1. 4. Compile a sample with a volume of 51,000 consonant phonemes.
 1. 5. Split the sample into 51 parts each comprising 1,000 phonemes.
 1. 6. Calculate the number of consonant phonemes for any position of the phoneme in a text.
 1. 7. Calculate the mean value of each consonant phoneme in a text with a volume of part that contains 1,000 phonemes and with a volume of the sample of 51,000 phonemes.
 1. 8. Combine consonant phonemes in groups (summing the mean frequencies of phonemes).

The result is the determined values of the mean frequencies of groups of consonant phonemes.

1. Check whether the mean frequencies of groups of consonant phonemes match the normal distribution by using the Pearson criterion:

1. 1. Determine a theoretically normal distribution.
1. 2. Calculate a theoretical frequency (the mathematical expectation that the magnitude of X is in the i -th interval).
1. 3. Check the conformity to the normal distribution for eight groups of phonemes (51 parts for each).

Provided that the mean frequencies of groups of consonant phonemes comply with the normal distribution, it is necessary to perform computerized style differentiation based on the mean groups of frequencies using the Student's criterion.

The algorithm of functioning of the system supports simultaneous work with two text files (Fig. 2). This includes opening two files, converting them into transcription, sampling of consonant phonemes, splitting the sample into portions, calculation of the number of phonemes in each portion and the sample, merging into groups and further verification

by the Pearson criterion. This is performed so that it is possible, provided the mean frequencies of groups of consonant phonemes comply with the normal distribution, to compare the texts for the existence of phonetical difference.

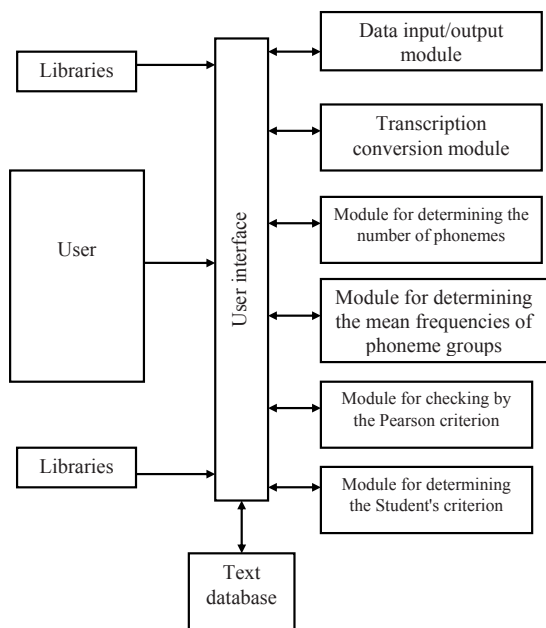


Fig. 1. Structure of software system for the differentiation of phonostatic structures of functional styles of the English language

In the process of software development we constructed the following basic classes: Main, Window, PanelFile, ExtFileFilter, PanelTranscription, DistributionOfPortion, DitributionOfGroup, CriterionPearson, CriterionStudent. The developed structure of classes enables choosing a text file, checking whether a given file has the .txt extension, con-

verting the text into a transcription variant. Input samples are checked by the system for conformity with the normal distribution law and are differentiated based on the mean frequencies of groups of consonant phonemes.

Using the java programming language ensures that the developed software is platform-independent.

7. Discussion of results of testing a system for the differentiation of phonostatic structures of styles

We have chosen as the material to study texts written in the literary, conversational, newspaper, and scientific styles. Specifically, Fig. 2 shows example of the interface for adding new words to the Word.txt and Transcription.txt files.

We tested the system using material of the texts written by different authors in the scientific style. In the “Pearson Criterion” tab we verified conformity of the texts to the law of normal distribution. It was established that groups of labial, front-alveolar, mid-alveolar, post-alveolar, nasal, sonorous, slit and closed phonemes comply with the law of normal distribution. Based on the differentiation of phonostatic structures of texts, by different authors, related to the scientific style, for the Student’s criterion, we established significant differences in styles for groups of labial, front-alveolar, post-alveolar, nasal, slit and closed phonemes. Random differences were found for groups of mid-alveolar and sonorous phonemes. Thus, we have established phonostatic parameters for the differentiation of texts by different authors.

Based on the research results, obtained for the scientific, fiction, conversational and newspaper styles, we determined significant substantial differences for the group of slit phonemes by the ranking method (rank indicators difference is 6). Fig. 3 shows statistical model of style differentiation for the scientific and conversational styles based on the ranking method for the group of slit phonemes for the case of an undefined position of the phoneme in a word:

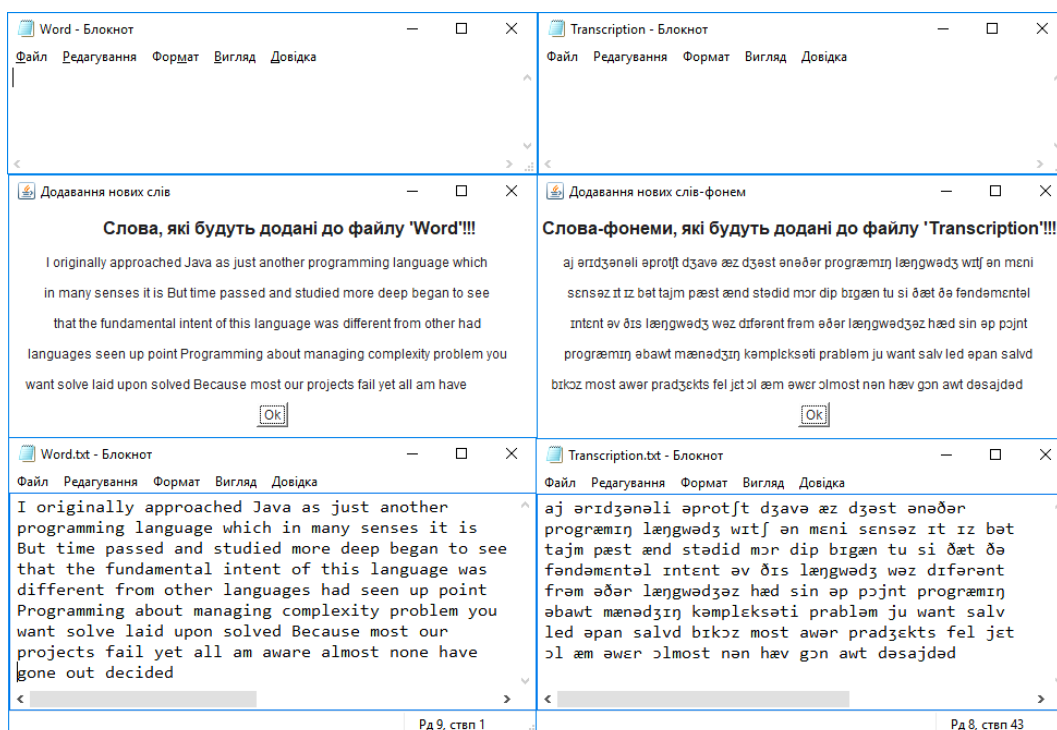


Fig. 2. Example of the interface for adding new words to the Word.txt and Transcription.txt files

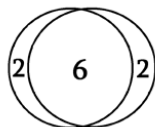


Fig. 3. Statistical model of style differentiation based on the ranking method: 6 – significant essential difference (6 units); 2 – insignificant similarity (2 units)

For the case of identifying the authorship of texts related to various subjects, but of one style and subtype, it is appropriate to apply a statistical model that combines three statistical models-elements: determining a style affiliation; determining a subtype affiliation; identifying an author of texts related to various topics. This is a statistical model for determining a general stylistic markedness of the examined text (Fig. 4).

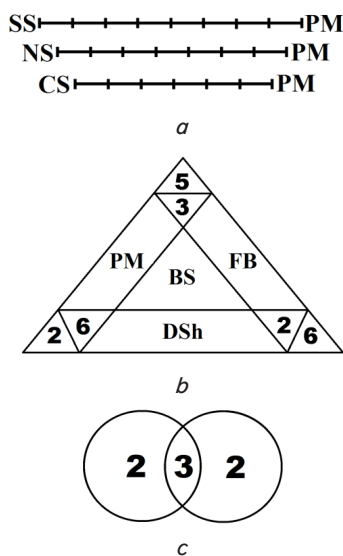


Fig. 4. Model: *a* – style differentiation for the case when a phoneme is at the beginning of a word when comparing texts of poems by Moore (PM), conversational (CS), newspaper (NS), and scientific styles (SS); *b* – subtype differentiation for the case when a phoneme is at the end of a word when comparing texts of poems by Byron and Moore, fiction by Byron (FB) and drama by Shaw (DSh); *c* – author’s differentiation for the case of the undefined position of the phoneme in a word when comparing texts of poetry by Byron and Moore; the belles letters style (BS)

The research results based on 5 out of 553 experiments (described earlier, in particular, in [22, 23, 26, 27]) showed that the developed methods, models, and tools make it possible to improve the efficiency of author’s attribution of a text. The phonological level selected for the study is organized stricter than the other levels of a language. However, the phonological system is probabilistic in character with the probability of making an error being equal to 5 %. The developed software system could be applied for identifying the authorship of a text in fiction, as well as legal, official, business, and scientific areas. Further study will address the

development of a software system for the author attribution of a text for each of the groups of consonant phonemes in order to determine a group of phonemes for which author attribution would be most effective.

8. Conclusions

1. Effectiveness of the methods developed was tested during 553 experiments, the results of five of which are covered in this paper. Experiments were conducted for eight groups of consonant phonemes (labial, front-alveolar, mid-alveolar, post-alveolar, nasal, sonorous, slit, and closed) in texts related to fiction, as well as conversational, newspaper, and scientific styles, for the three cases of the position of a phoneme in a word. The results of experiments, described previously and examined in present work, show that the developed method of a comprehensive analysis of the differentiation of phonostatic structures of styles, as well as a multi-factor method for determining the degree of action of factors related to style, subtype, and the author’s manner of presentation, make it possible to improve the efficiency of author attribution, and thereby check a text for plagiarism. Efficiency of the method for a comprehensive analysis of differentiation of the phonostatic structures of styles is ensured by the proposed combination of statistical methods (hypotheses, ranking, determining distances between styles), among which the ranking method was applied for the first time to solve the task on author attribution of a text. Data were obtained from three methods of mathematical statistics with a probability of error of 5 %. Efficient is the scheme style→subtype→author, which underlies a multifactor method for determining the degree of action of factors related to style, subtype, and an author’s manner of presentation.

2. Based on the developed methods, we built a statistical model of style differentiation using the ranking method and a statistical model for determining a general stylistic markedness of the examined text. The models allow the improvement of accuracy in the differentiation of phonostatic structures of styles during author attribution and verifying a text for plagiarism.

3. We have developed the structure and tools for a system of differentiation of phonostatic structures of styles, which is different from existing ones by the chosen level of a language – phonological. At this level of a language one can obtain results with a greater accuracy. The constructed system is based on a modular principle, which makes it possible to rapidly modify the developed software and to identify a group of consonant phonemes, which could be employed to perform author attribution of a text more effectively. The system was implemented in the programming language java, which ensures that it is platform-independent. The developed and implemented system can operate at different computer platforms.

The research results obtained could be used for identifying the authorship of the examined text, as well as for verifying a text for plagiarism. Further research seems promising in terms of defining phonostatic parameters, specifically, the style differentiating power of groups of consonant phonemes whose mean frequencies are the criterion for the differentiation of an author’s style.

References

1. Kornai A. Mathematical Linguistics. Springer, 2008. doi: 10.1007/978-1-84628-986-6
2. Gries Th. S. Statistics for Linguistics with R. Mouton Textbook, 2009. 335 p. doi: 10.1515/9783110216042

3. Martindale C., McKenzie D. On the utility of content analysis in author attribution: The Federalist // *Computers and the Humanities*. 1995. Vol. 29, Issue 4. P. 259–270. doi: 10.1007/bf01830395
4. Gibbons J. *Forensic Linguistics. An Introduction to Language in the Justice System*. Wiley-Blackwell, 2003. 346 p.
5. Olsson J. *Forensic Linguistics. Second edition: An Introduction to Language, Crime and the Law*. Bloomsbury Academic, 2008. 288 p.
6. Berko A. Yu., Vysotska V. A., Chyrun L. V. Linhvistychnyi analiz tekstovoho komertsynoho kontentu // *Informatsiyni systemy ta merezhi. Visnyk Natsionalnoho universytetu "Lvivska politekhniky"*. 2015. Issue 814. P. 203–227.
7. Bisikalo O. V., Vysotska V. A. Sentence syntactic analysis application to keywords identification Ukrainian texts // *Radio Electronics, Computer Science, Control*. 2016. Issue 3. P. 54–65. doi: 10.15588/1607-3274-2016-3-7
8. Shakhovska N., Vysotska V., Chyrun L. Intelligent Systems Design of Distance Learning Realization for Modern Youth Promotion and Involvement in Independent Scientific Researches // *Advances in Intelligent Systems and Computing*. 2016. P. 175–198. doi: 10.1007/978-3-319-45991-2_12
9. Content linguistic analysis methods for textual documents classification / Lytvyn V., Vysotska V., Veres O., Rishnyak I., Rishnyak H. // 2016 XIth International Scientific and Technical Conference Computer Sciences and Information Technologies (CSIT). 2016. doi: 10.1109/stc-csit.2016.7589903
10. Lytvyn V. V., Bobyk I. O., Vysotska V. A. Application of algorithmic algebra system for grammatical analysis of symbolic computation expressions of propositional logic // *Radio Electronics, Computer Science, Control*. 2016. Issue 4. P. 77–89. doi: 10.15588/1607-3274-2016-4-10
11. Development of a method for the recognition of author's style in the Ukrainian language texts based on linguometry, stylemetry and glottochronology / Lytvyn V., Vysotska V., Pukach P., Bobyk I., Uhryn D. // *Eastern-European Journal of Enterprise Technologies*. 2017. Vol. 4, Issue 2 (88). P. 10–19. doi: 10.15587/1729-4061.2017.107512
12. Davydov M., Lozynska O. Linguistic models of assistive computer technologies for cognition and communication // 2016 XIth International Scientific and Technical Conference Computer Sciences and Information Technologies (CSIT). 2016. doi: 10.1109/stc-csit.2016.7589898
13. Modelling of semantics of natural language sentences using generative grammars / Shestakevych T., Vysotska V., Chyrun L., Chyrun L. // *Computer Science and Information Technologies: Proc. of the IX-th Int. Conf. CSIT'2014*. Lviv: Lviv Polytechnic Publishing House, 2014. P. 19–22.
14. Application of sentence parsing for determining keywords in Ukrainian texts / Vasyl L., Victoria V., Dmytro D., Roman H., Zoriana R. // 2017 12th International Scientific and Technical Conference on Computer Sciences and Information Technologies (CSIT). 2017. doi: 10.1109/stc-csit.2017.8098797
15. Zhezhnych P., Markiv O. A linguistic method of web-site content comparison with tourism documentation objects // 2017 12th International Scientific and Technical Conference on Computer Sciences and Information Technologies (CSIT). 2017. doi: 10.1109/stc-csit.2017.8098800
16. Peculiarities of content forming and analysis in internet newspaper covering music news / Korobchinsky M., Chyrun L., Chyrun L., Vysotska V. // 2017 12th International Scientific and Technical Conference on Computer Sciences and Information Technologies (CSIT). 2017. doi: 10.1109/stc-csit.2017.8098735
17. Kapociute-Dzikiene J., Utkia E., Sarkute L. Authorship Attribution and Author Profiling of Lithuanian Literary Texts // *Proceedings of the 5th Workshop on Balto-Slavic Natural Language Processing*. Hissac, Bulgaria, 2015. P. 96–105.
18. Stamatatos E. A survey of modern authorship attribution methods // *Journal of the American Society for Information Science and Technology*. 2009. Vol. 60, Issue 3. P. 538–556. doi: 10.1002/asi.21001
19. Automatically profiling the author of an anonymous text / Argamon S., Koppel M., Pennebaker J. W., Schler J. // *Communications of the ACM*. 2009. Vol. 52, Issue 2. P. 119. doi: 10.1145/1461928.1461959
20. Koppel M., Schler J., Argamon S. Computational methods in authorship attribution // *Journal of the American Society for Information Science and Technology*. 2009. Vol. 60, Issue 1. P. 9–26. doi: 10.1002/asi.20961
21. Juola P. Authorship Attribution // *Foundations and Trends® in Information Retrieval*. 2007. Vol. 1, Issue 3. P. 233–334. doi: 10.1561/1500000005
22. Khomytska I., Teslyuk V. The Method of Statistical Analysis of the Scientific, Colloquial, Belles-Lettres and Newspaper Styles on the Phonological Level // *Advances in Intelligent Systems and Computing*. 2016. P. 149–163. doi: 10.1007/978-3-319-45991-2_10
23. Khomytska I., Teslyuk V. Specifics of phonostatistical structure of the scientific style in English style system // 2016 XIth International Scientific and Technical Conference Computer Sciences and Information Technologies (CSIT). 2016. doi: 10.1109/stc-csit.2016.7589887
24. Bektaev K. B. *Matematicheskie metody v yazykoznanii*. Ch. 2. Alma-Ata, 1974. 335 p.
25. Mitropol'skiy A. K. *Tekhnika statisticheskikh vychisleniy*. Moscow: Nauka, 1971. 576 p.
26. Khomytska I., Teslyuk V. Modelling of phonostatistical structures of English backlingual phoneme group in style system // 2017 14th International Conference The Experience of Designing and Application of CAD Systems in Microelectronics (CADSM). 2017. doi: 10.1109/cadsm.2017.7916144
27. Khomytska I., Teslyuk V. Modelling of phonostatistical structures of the colloquial and newspaper styles in english sonorant phoneme group // 2017 12th International Scientific and Technical Conference on Computer Sciences and Information Technologies (CSIT). 2017. doi: 10.1109/stc-csit.2017.8098738
28. Chabanyuk Y., Seniv M., Khimka U. Continuous Stochastic Optimization Procedure in Software Reliability // *Proceedings of the XIIth International Conference The Experience of Designing and Application of CAD Systems in Microelectronics CADSM 2013*. Polyana, 2013. P. 56–59.