# HYBRID METHOD OF SPEECH SIGNALS SPECTRAL ANALYSIS BASED ON THE AUTOREGRESSIVE MODEL AND SCHUSTER PERIODOGRAM

**V. V. Savchenko**                                    UDC 53.082.4; 519.246.87; 004.934.2

*The task of measuring the spectral density of power of a speech signal in sliding observation window mode is examined. A parametric approach to solving this task using an autoregressive data model is studied. The problem of optimizing the order of an autoregressive model under the conditions of small samples is studied. It is proposed to solve the problem using a hybrid method of spectral analysis based on sequential enumeration of a finite number of variants. The optimization criterion is formulated in terms of an inverse problem: from the speech signal to the voice source. It uses the scale-invariant measure of the spectral distance as the objective function, and the Schuster periodogram as the reference sample. The effectiveness of the hybrid method has been experimentally evaluated on the basis of the author's software. It is shown that with the duration of the observation window no greater than 10 ms, the use of the hybrid method increases the accuracy of spectral analysis by more than 30%, compared to the well-known Berg method, the order of which is established according to the Akaike information criterion.*

***Keywords:*** *acoustic measurements, speech acoustics, speech signal, autoregressive model, small samples, a priori uncertainty, adaptive approach.*

**Introduction.** Spectral analysis of speech signals is related to to the most dynamically developing directions of study in the field of acoustic measurements. With reference to digital spectral analysis, the problem is, as a rule, reduced to determining the envelope of a power spectrum in the interval of quasi-stationarity of a speech signal of duration one frame [1, 2]. The spectral envelope is responsible for the timbre and voice recognition of the speaker [3–5]. Here the small samples problem arises [6–8]. The theory recommends parametric methods of spectral analysis to overcome this, such as the Berg method, autocorrective, covariational, and so forth. [9, 10] These may be studied as variants to classical methods [11] based on the discrete Fourier transform. The autoregressive model [4] used in the parametric methods corresponds well with the "acoustic pipe" model of the vocal pathway of the speaker [2]. However, application of the autoregression model does not fully solve the small sample problem, which is merely transformed into another acute problem of the parametric methods of spectral analysis: the optimization of order $p < \infty$ of the autoregressive model [3, 12]. This optimization is usually associated with the Akaike informational measure and analogs [13–15]. However, when processing speech signals in the observational sliding window mode of small duration $\tau = 5$–$30$ ms, the effectiveness of these criteria sharply decreases [16, 17].

The objective of this article is to study the capabilities and effectiveness of a hybrid technique under conditions of small samples of observations. The subject of the study is a hybrid method of the spectral analysis of speech signals which unites the advantages of non-linear parametric and classical linear methods. The problem is formulated as optimization in the terms of the multivariant verification of $R$ hypotheses $p = p_\mathrm{r}$, $r = \overline{1,\,R}$ on the order $p \geq 1$ of the autoregressive model. The criterion will be the principle of the maximum level of correspondence of the autoregressive model of the Schuster periodogram as the most informative characteristic relative to the fine structure of a speech signal [3].

This article is written for development of the results of previous work of the author [4, 10], performed in co-authorship with the staff of the Laboratory of Algorithms and Techniques for the Analysis of Network Structures of NIU "Higher School of Economics".

**Problem definition.** We consider a speech signal $\{x(n)\}$ in the discrete time $n = 0,\,1,\,\ldots,\,N-1$ in the observation interval of finite duration $\tau < \infty$, where $N = \tau F$ is the sample size; and $F$ is the frequency of time sampling. We define the spectral density of power (SDP) of a signal through an autoregressive model of the general form [1]:

$$G_p\left(f\right) = \sigma_p^2 T \left|1 - \sum_{i=1}^{p} a_p\left(i\right) \exp\left(-\mathrm{j}2\pi i f T\right)\right|^{-2}, \tag{1}$$

where $a_p(i)$ is the $i$th element of the $p$-vector $\mathbf{a}_p$ of the coefficients of the linear autoregression of the $p$th order;     is the scale factor; and $T = F^{-1}$.

The parameters $\mathbf{a}_p$ and $\sigma_p^2$ are defined in (1) from the sample $\{x(n)\}$ by means of one of the known methods of parametric spectral analysis, such as the Berg method, autocorrective, covariational, and so forth [1]. However, in any variant the order $p$ must be set a priori. The problem consists in the fact that depending on the value of $p$ established in (1), the form of the autoregressive model of a speech signal varies, and especially strongly in the conditions of small samples [3]. Taking this into account, the problem is reduced to determining the measure of optimization of the order of a spectral estimate (1) in the sliding window mode for observations of finite duration $\tau$.

**Criterion of optimization.** Since an autoregressive model (1) is inseparably linked with the envelope of the SDP of a speech signal [17], we will define the pectrum of signal power through the discrete Fourier transform [2]:

$$G_x\left(f_m\right) = \frac{\left|\sum_{n=0}^{M-1} x\left(n\right) \exp\left(-\mathrm{j}2\pi n m M^{-1}\right)\right|^2}{FM}, \quad m = 0, 1, \ldots, M - 1, \tag{2}$$

where $f_m = m\,\Delta f$ is the discrete frequency with a shift of $\Delta f = FM^{-1} = \mathrm{const}$ in relation to the frequency $f_{m-1}$; and $M$ is the size of the set of readings of a speech signal in the frequency domain ($M < \infty$).

In this theory, expression (2) is known as a Schuster periodogram [1]. This periodogram of dimensionality $M = 2^k \geq N$, where $k$ is some integral number, is calculated using algorithms for the fast Fourier transform [18]. Here only the first $N$ from $M$ ($N \ll M$) readings of the speech signal in expression (2) are distinct from zero. For example, for a duration of the observation window $\tau = 5$–$30$ ms and sampling frequency $F = 8$ kHz[1] for the case $k = 10$, we will have $M = 1024$ with the equality $N = 40$–$240$ (this is the typical relation of the number of readings of a speech signal in the frequency and time fields, respectively) [10].

Due to the linear nature of the discrete Fourier transform, expression (2) contains the maximal useful information about the SDP, and therefore it can be used as a spectral reference sampler in the problem of optimizing the autoregressive model. For this purpose, we rewrite (1) in the form of an inverse conversion

$$G_p\left(f_m\right) = \sigma_p^2 T \left|A_p\left(\mathrm{j}f_m\right)\right|^{-2}$$

of the square of the modulus of the complex coefficient of transmission

$$A_p\left(\mathrm{j}f_m\right) = 1 - \sum_{i=1}^{p} a_p\left(i\right) \exp\left(-\mathrm{j}2\pi i m M^{-1}\right), \quad m = \overline{0;\ M-1}, \tag{3}$$

of the transversal filter of the $p$th order on the extracted sample of discrete frequency $\{f_m\}$. The signal $\{y(n)\}$ at the exit of such a filter in the frequency domain is described by expression [4]:

$$G_y\left(f_m; p\right) = \left|A_p\left(\mathrm{j}f_m\right)\right|^2 G_x\left(f_m\right) = \sigma_p^2 T \frac{G_x\left(f_m\right)}{G_p\left(f_m\right)}, \quad m = \overline{0; M-1}. \tag{4}$$

---

[1]In accordance with the pass band of a standard telephone communications channel.

Expression (4) defines the operation of smoothing or whitewashing the envelope $\overline{G}_x(f_m)$ of the power spectrum at input (2). Here, in the ideal for the spectral envelope of $\overline{G}_y(f_m;p)$ at exit of the equalizer network (3), the system of equations $\forall m < M$: $\overline{G}_y(f_m; p) = G_0 = \mathrm{const}$ is valid. In terms of the inverse problem [3] "from a speech signal to its voice radiant," this defines the hypothetical white noise as the generating process for the autoregressive model (1).

However, in practice the form of the envelope of the SDP (4) may differ substantially from rectangular, due first of all to the non-ideal nature of the autoregressive model used. (1). The basic value in this sense has order $p$. We optimize its value within the finite set ($R$-set) of the variants $p_1 < p_2 < ... < p_R$ by the principle of the maximum likelihood of the form of the envelope of the SDP (4) to rectangular.

Taking into account the above-specified, we use the modified standard of the COSH distance (COSH is the hyperbolic cosine) function [19] as the function of the purpose of the optimization problem being solved:

$$\rho\,(p) = \sqrt{\left[M^{-1}\sum_{m=0}^{M-1} G_0\overline{G}_y^{-1}\,(f_m;p)\right]\left[M^{-1}\sum_{m=0}^{M-1}\overline{G}_y\,(f_m;p)\,G_0^{-1}\right]} - 1$$

$$= M^{-1}\sqrt{\left[\sum_{m=0}^{M-1}\overline{G}_y^{-1}\,(f_m;p)\right]\left[\sum_{m=0}^{M-1}\overline{G}_y\,(f_m;p)\right]} - 1 \geq 0\,,$$

(5)

with the property of scale invariance to the SDP of a speech signal at input. With equality of the envelope $\overline{G}_y(f_m;p)$ to an arbitrary constant $G_0$, measure (5) will be identically equal to zero. According to the results of [3], the limit accuracy of the autoregressive model (1) will be reached in this case. As measure (5) increases, the accuracy of the autoregressive model decreases, and therefore measure (5) may be examined as an objective index of the error of the generated autoregressive model. From this follows the criterion for adopting solutions in the problem of spectral analysis of a speech signal: on the set $R$ of examined variants $\{p_r, r \leq R\}$, the optimal estimate of the SDP $\{G^*(f_m)\}$ corresponds to expression (1) in the selection of the order of an autoregressive model according to the rule

$$p^* = \mathrm{Arg}\min_{\forall r \leq R}\ \rho\,(p)\big|_{p=p_r}\,.$$

(6)

Practical implementation of the proposed criterion reduces mainly to an estimate of the envelope of the SDP at exit of the leveling network (3) and to calculation of the modified COSH distance (5) within the set of variants $\overline{G}_y(f_m;\,p_r)$, $r = \overline{1,\,R}$, in the interval of observations of finite duration $\tau$. Here, is necessary to use fast computational methods.

**Example of practical implementation.** Turning away from the Berg method [20], which possesses at the same time fast response and high resolution capability by frequency, we use the Levinson recursion [21] to adapt the vector of coefficients of the SDP (1) for the sample $\{x(n)\}$ of finite size $N$:

$$\forall q = \overline{1,\,p}:\ a_q\,(i) = a_{q-1}\,(i) + c_q a_{q-1}\,(q-i)\,,\ i = \overline{1,\,q}\,;$$

(7)

$$c_q = \frac{2\sum\limits_{n=q}^{N-1}\eta_{q-1}\,(n)\,v_{q-1}\,(n-1)}{\sum\limits_{n=q}^{N-1}\left[\eta_{q-1}^2\,(n) + v_{q-1}^2\,(n-1)\right]}\,;$$

$$\eta_q\,(n) = \eta_{q-1}\,(n) + c_q v_{q-1}\,(n-1)\ ;$$

$$v_q\,(n) = v_{q-1}\,(n-1) + c_q \eta_{q-1}\,(n)\ ;$$

$$v_0\,(n) = \eta_0\,(n-1) = x\,(n)\,,\ n = 0, 1, ..., N-1\,.$$

Here, in accordance with [1], the stability of the formulated autoregressive model, in the sense of digital filtering [22] independently of the sampling composition, is guaranteed. The scale factor $\sigma_p^2$ does not play a large role in expression (1). As a rule [4], this is established from the condition of normalizing the estimate of the SDP on variance to the level set by the user.

For selection of the spectral envelope (4) by analogy with [3], the standard recirculator is used. Its dynamic in the frequency domain is described by the difference equation [18]:

$$\overline{G}_y \left( f_m; \, p \right) = b\overline{G}_y \left( f_m - \Delta f; \, p \right) + G_y \left( f_m; \, p \right), \quad m = \overline{1, M} , \tag{8}$$

where $b$ = const.

The operation of the recirculator reduces to accumulating the spectral components of the speech signal in sliding window mode in the frequency domain. The size of a window defines the inertance of the reciirculator which is governed in (8) by the parametric value of $b$. Depending on the type of analyzed speech frame (vocalized[2] or not), the quantitative index of inertance $\theta = -\Delta f/\ln b = -FM^{-1}/\ln b$ must be in agreement either with the frequency of the fundamental component of the speech signal [2], or with the resolution capability of the spectral analyser at frequency $\theta = \tau^{-1} = FN^{-1}$. The first variant was described in sufficient detail in [3]. For this reason, we will later examine the spectral analysis of non-vocal frames containing voiceless consonants of the speech of a conventional speaker. For this variant [4]:

$$b = \exp \left( -\tau FM^{-1} \right) = \exp \left( -NM^{-1} \right). \tag{9}$$

For example, for $\tau$ = 5–30 ms, $F$ = 8 kHz, and $M$ = 1024, we obtain $b$ = 0.80– 0.96.

Expressions (1)–(9) in the aggregate define the hybrid technique of the spectral analysis of a speech signal $(x (n)\}$ using at the same time the Schuster periodogram and the Berg method as a methodological basis. The proposed hybrid technique was for the first time implemented in practice, based on author software,[3] using which the experiment described later was set up and conducted.

**Program and results of experiment.** The object of the experimental study was a synthesized 10th-order autoregressive process, imitating the consonant "sh"of the Russian alphabet in the speech of the speaker being studied (the author of this article) and specified by the vector of autoregression coefficients $\mathbf{a}_{10}^*$ = (1.708427645; 2.186231624; 2.390214916; 2.155074921; 1.246001575; 0.49162809; –0.068969197; –0.416448837; –0.334798662, –0.08981284). The generating process for the autoregression process was white Gaussian noise with variance $\sigma_{10}^2$. The synthesized signal of sufficient duration $T_0$ = 3 s was linearly divided in the Phoneme Training program into frames $\{x(n)\}$ of duration $\tau$, equal to 5, 10, and 30 ms. As a result, for each of the three durations $\tau$ of the observation window, a representative database up to $L = T_0/\tau$ = 300 independent frames was formulated. For a signal sampling frequency $F$ = 8 kHz, the sample size $N$ for each such frame was 40, 80, and 240 readings, respectively. Further, for each frame, the Schuster periodogram (2) of dimensionality $M$ = 1024 and frequency step $\Delta f$ = 8000/1024 = 7.8125 Hz was calculated using the fast Fourier transform. At the same time with the Schuster periodogram, the Berg method obtained a spectral estimate (1) at $R$ = 16 alternate variants of its order $p$ = 5–20. The total computational complexity of the computational procedure (7) was on the order of $3Np_R = 3 \cdot 240 \cdot 20 = 14.4 \cdot 10^3$ elementary operations in the interval of duration $\tau$. This, with a margin, is responsible for the efficiency of modern information systems and real-time technologies [18, 19]. Further, according to expressions (8) and (9), estimates are obtained for the spectral envelope $\overline{G}_y(f_m; \, p_r)$ of the signal $\{y(n)\}$ at output of the leveling filter (3), for all variants $p_r, r \le R$. The corresponding value of the characteristics of $\rho(p)$ is determined from these estimates in accordance with (6). The obtained values were averaged further from the results of $L$ independent tests of observations for each duration $\tau$. As a result the relative error of the experimental data [4, 5] was no greater than $\varepsilon = 165 \cdot (300)^{-1/2}$ = 9.5% with a confidence probability 0.9 and above.

Figure 1 shows a typical family of Berg spectral estimates obtained in the same observation window of duration $\tau$ = 10 ms. The corresponding Schuster periodogram (2) is shown by a dashed line for comparison One may observe the advantages of the variant of a spectral estimate of order $p$ = 10. They are confirmed by the graph of the corresponding SDP (4) in Fig. 2, where the dashed line reflects the form of the spectral envelope: for $p$ = 10, the form is maximally close to

---

[2]Vocalized frames are differentiated and selected from the speech stream according to the amplitude indication [4, 10].

[3]Phoneme Training. URL: https://sites.google.com / site / frompldcreators / produkty-1 / phonemetraining (date: accessed 2/2/2023).
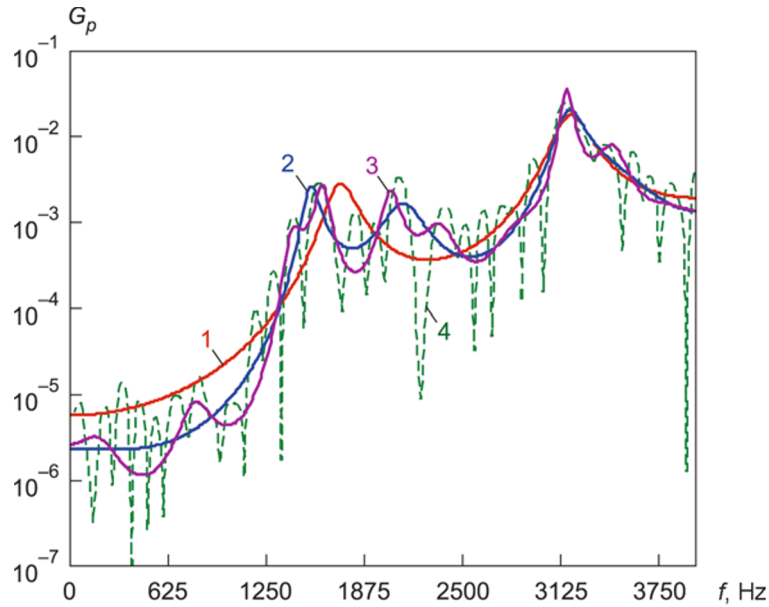
Fig. 1. A family of Berg spectral estimates of order $p = 5$, 10, and 20 (curves 1–3) against the background of the Schuster periodogram (curve 4) at duration $\tau = 10$ ms.
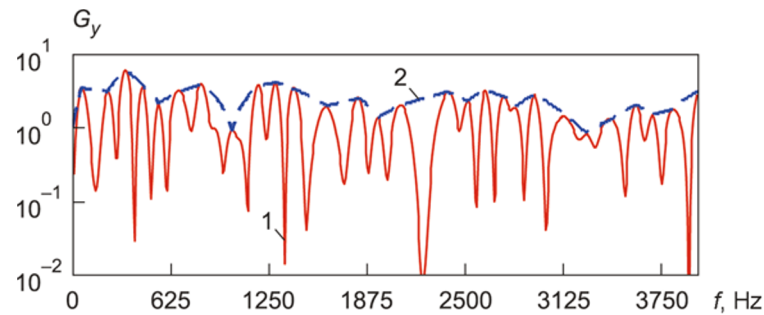


Fig. 2. Spectral power density of the signal at the output of the leveling filter (4) (curve 1) of order $p = 10$ at duration $\tau = 10$ ms and its spectral envelope (curve 2).

rectangular. The histogram in Fig. 3 — the operating characteristics (5) of the hybrid techniques for $\tau = 10$ ms — serves as rigorous substantiation of order $p = 10$. In accordance with the histogram, a strong argument for the use of criterion (6) is the global minimum $\rho$ (10) ≈ 0.085 of the operating characteristicss $\rho(p)$. Under the examined conditions, the minimum point also defines the best value $p^* = 10$ of the order of the autoregressive model (1).

The conclusions drawn are not trivial, if one comparea the behavior of the two operating characteristics: the proposed criterion $\rho(p)$ and the Akaike information criterion in its modified variant of BIC (the Bayesian information criterion) [14]:
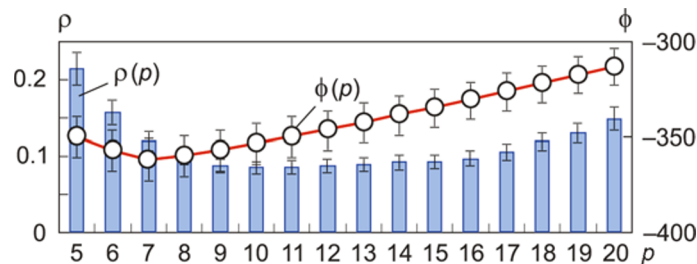


Fig. 3. Working characteristics of the criterion $\rho(p)$ (6) and {BIC} $\phi(p)$ at duration $\tau = 10$ ms.

207

$$\phi\left(p\right) \triangleq N \ln \sigma_p^2 + p \ln N .$$

For $\tau = 10$ ms, the dependence $\phi(\rho)$ is also provided as a scatter plot in Fig. 3, and the minimum of the dependence is reached at the point $p = 7$. This is a substantially different result in comparison with the order established above, $p^* = 10$. As a result, in the case being studied, one may speak of the gain in precision of the hybrid technique in comparison with the Berg method when applying the Akaike criterion, defined as [3]:

$$B\left(p\right) = \frac{\rho\left(p\right) - \rho\left(p^*\right)}{\rho\left(p\right)} \cdot 100 = \frac{\rho\left(7\right) - \rho\left(10\right)}{\rho\left(7\right)} \cdot 100 = \frac{0.12 - 0.085}{0.12} \cdot 100 \simeq 29.2\% .$$

When the duration of a speech frame is reduced to $\tau = 5$ ms, the gain increases to $B(p) = 36.6\%$. In order to illustrate the obtained results, Fig. 4a shows the family of operating characteristics (histograms) $\rho(p)$ of the hybrid technique for three observational durations. These characteristics are similar in nature, and a tendency is observed towards smoothing the base in neighborhoods of the minimum point $p^* = 10$ s, as the duration $\tau$ of a frame increases. This fact is explained by the known universality of higher-order autoregressive models [1].
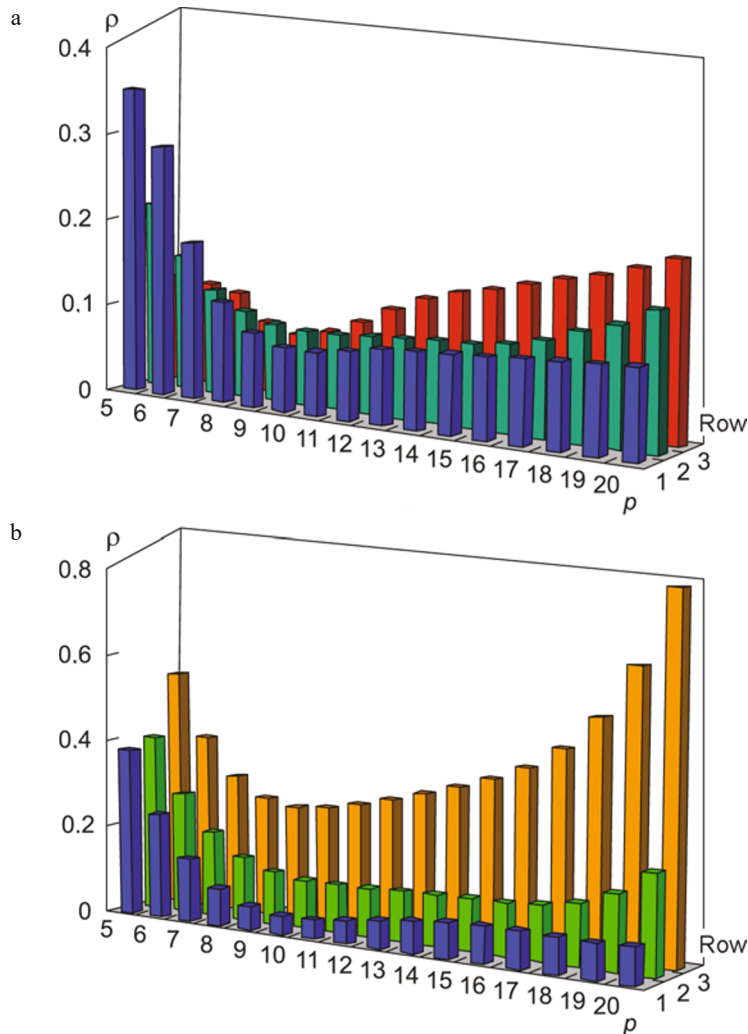


Fig. 4. A family of histograms of the operating characteristics (6) (a) and the dependencies (10) (b) at $\tau = 30$, 10, and 5 ms (rows 1–3, respectively).
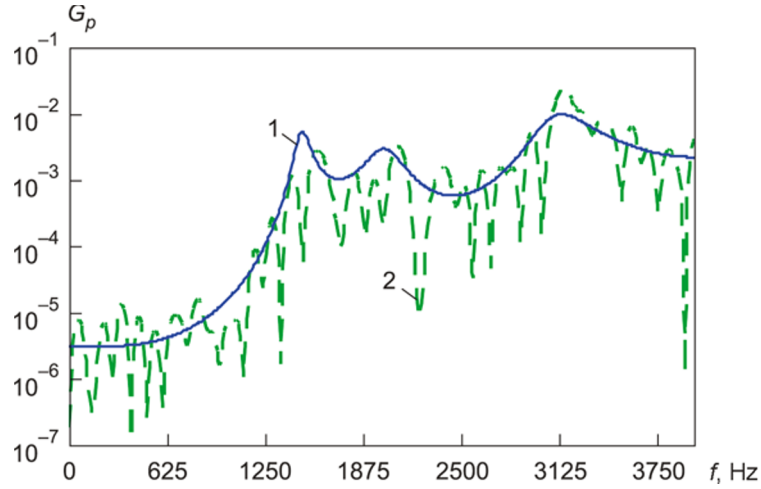
Fig. 5. The true power spectral density of the signal at the smoothing filter at input (3) (curve 1) against the background of the Schuster periodogram (curve 2) at $\tau = 10$ ms.

**Discussion of obtained results.** The results of the second and final stage of the experimental studies can serve as additional reasons for the proposed hybrid technique of spectral analysis and its criterion (6). At this stage, the autoregressive model (1) was adapted to sample the observations by the Berg method in the same condition as in the first stage. In variant (3), the autoregressive model was also used in order to smooth the envelop of the Schuster periodogram (2). However the envelope of the SDP $G_y(f_m; p)$ was compared, according to its form in (6), not with the rectangular envelope $G_0 = $ const, but with the envelope $\overline{G}_y(f_m; p^*)$ of the SDP at exit of the ideal smoothing filter (3) for which the vectorial equality $\mathbf{a}_p = \mathbf{a}_{10}^*$ is fulfilled. Figure 5 presents a graph of the true SDP (1) against the same Schuster periodogram that is portrayed in Fig. 1. By analogy with (5), we write as the characteristics in this variant

$$\rho\left(p\right) = M^{-1} \sqrt{\left[ \sum_{m=0}^{M-1} \overline{G}_y^*\left(f_m;\, p\right) \overline{G}_y^{-1}\left(f_m;\, p\right) \right]} \sqrt{\left[ \sum_{m=0}^{M-1} \overline{G}_y\left(f_m;\, p\right) \overline{G}_y^{*-1}\left(f_m;\, p\right) \right]} - 1\,. \tag{10}$$

This is only a hypothetical variant, unrealizable in practice due to the definition of the spectral analysis problem under conditions of *a priori* uncertainty. There is interest in the comparison of dependence (10) with the characteristics (5) within the criterion used (6).

Figure 4b shows the family of histograms of the dependence of $\rho(p)$ for the three durations of observations that were studied, from which follows the practically perfect analogy with the data of Fig. 4a, Hence, one may derive a conclusion on successful resolution of the small samples problem by using the proposed hybrid technique.

**Conclusion.** In the proposed hybrid technique of the spectral analysis of speech signals in the sliding window of observations mode, the Berg spectral estimate interacts with the Schuster periodogram estimate, although in practice they usually compete with each other. As a result, it was possible, in the hybrid technique, to consolidate the known advantages of the Berg method (speed of convergence and resolution capability by frequency) and Schuster (precision of representations of the thin structure of a speech signal in the frequency domain). The results of the experiment that was performed serve as proof.

The results obtained may be used in development of systems of digital spectral analysis of speech signals, as well as of signals of a speech-like structure for the fields of technical, economic, and biomedical diagnostic [6–9, 12].

**Conflict of interest**.The author declare no conflict of interest.

# REFERENCES

1. S. L. Marple Jr. *Digital Spectral Analysis*. 2nd ed., Dover Publications, New York (2019).
2. L. R. Rabiner and R. W. Shafer, *Theory and Applications of Digital Speech Processing,* Pearson, Boston (2010).

3. V. V. Savchenko, Perfecting the Methodology of Measuring the Index of Accuracy of an Autoregressive Model of a Speech Signal, *Izmer. Tekh.,* No. 10, 58–63 (2022), https://doi.org/10.32446/0368-1025it.2022-10-58-63.

4. A. V. Savchenko and V. V. Savchenko, An Adaptive Method of Measuring the Frequency of the Fundamental Tone Using a Two-Level Autoregressive Model of a Speech Signal, *Izmer. Tekh.*, No. 6, 60–66 (2022), https://doi.org/10.32446/0368-1025it.2022-6-60-66.

5. V. V. Savchenko, *Radioelectron. Commun. Syst.*, **63**, 532–542 (2020), https://doi.org/10.3103/S0735272720100039.

6. Sh. Ando, *J. Acoust. Soc. Am.*, **146**, 2846 (2019), https://doi.Org/10.1121/1.5136873.

7. Yu. Gu and H. L. Wei, *Inform. Sci.*, **451–452**, 195–209 (2018), https://doi.Org/10.1016/j.ins.2018.04.007.

8. C. A. Liu, B.-S. Kuo, and W. J. Tsay, Autoregressive Spectral Averaging Estimator, *IEAS Working Paper*, No. 17-A013, available at: <https://www.econ.sinica.edu.tw /~econ/ pdfPaper/17-A013.pdf> (accessed: 2/2/2023) (2017).

9. A. A. Kuznetsov, Structural Frequency Analysis of Rhythmograms of Ill People, *Izmer. Tekh.,* No. 4, 46–51 (2014).

10. V. V. Savchenko and L. V. Savchenko, A method of Autoregressive Modeling of a Speech Signal Based on the Discrete Fourier Transform and Scale-Invariant Measurements of Informational Error, *Radiotekh. Élektron.,* **66**, No. 11, 1100–1108 (2021), https://doi.org/10.31857/S0033849421110085.

11. T. C. Mills, Schuster, Beveridge, and Periodogram Analysis, in *The Foundations of Modern Time Series Analysis,* pp. 18–29, Macmillan, London (2011), https://doi.org/10.1057/9780230305021_3.

12. A. V. Kashin, N. S. Kornev, N. A. Makarichev, et al., *Instrum. Exp. Tech.+*, **63**, 34–40 (2020), https://doi.org/10.1134/S0020441220010030.

13. V. V. Savchenko and A. V. Savchenko, *Radioelectron. Commun. Sys.,* **62**, No. 5, 223–231 (2019), https://doi.org/10.3103/S0735272719050042.

14. J. Ding, V. Tarokh, and Y. Yang, *IEEE Trans. Inform. Theory*, **64**, No. 6, 4024–4043 (2018), https://doi.org/10.1109/TIT.2017.2717599.

15. A. Boisbunon, S. Can, D. Fourdrinier, W. Strawderman, and M. T. Wells, *Int. Stat. Rev.*, **82**, No. 3, 422–439 (2014), https://doi.org/10.1111/insr.12052.

16. V. V. Savchenko, *Radioelectron. Commun. Sys.*, **64**, 592–603 (2021), https://doi.Org/10.3103/S0735272721110030.

17. M. Tohyama, Spectral Envelope and Source Signature Analysis, in *Acoustic Signals and Hearing,* pp. 89–110, Academic Press, (2020), https://doi.org/10.1016/B978-0-12-816391-7.00013-9.

18. Y. D. Shirman ed., *Radioelectronic Systems. Foundations of Construction and Theory: Reference Manual*, Radiotekhnika, Moscow (2007).

19. A. V. Savchenko, V. V. Savchenko, and L. V. Savchenko, *Optim. Lett.*, No. 16, 2095–2113 (2022), https://doi.org/10.1007/s11590-021-01790-5.

20. J. P. Burg, A New Analysis Technique for Time Series Data, in *Modern Spectral Analysis,* IEEE Press, New York (1978).

21. D. Xiao, E. Mo, Ya. Zhang, M. Zhao, and L. Ma, *Heliyon*, **4**, No. 11, Article ID e00948 (2018), https://doi.Org/10.1016/j.heliyon.2018.e00948.