# Problem of Reliability Justification of Computation Error Estimates

## V. P. Zhitnikov[1]

## N. M. Sherykhalina[2]

## A. A. Sokolova[3]

[1] *Ufa State Aviation Technical University, Dr, Professor*
[2] *Ufa State Aviation Technical University, Dr, Professor*
[3] *Ufa State Aviation Technical University, Postgraduate Student*
*Adress: 450000, K. Marx str., 12, Ufa, Russia*
*E-mails: zhitnik@mail.ru; n_sher@mail.ru, inform-rutaxi.mail.ru*

*Abstract*

*The problem of reliability justification of errors estimates of results obtained by numerical methods is under consideration. The application possibility conditions of the known methods and their downsides are discussed. A new method is offered including numerical filtration of results, choice of a standard and correctness verification of its choice. This method allows the obtainment and visualization of certain error estimates and justification of result reliability. The numerical solutions obtained by different methods of the Stokes solitary wave problem are shown. Also, the given problem results computed by different authors with help of the known methods are cited. Computations possibility of results with high accuracy and preference of the suggested approach for estimate reliability justification are demonstrated by the sample of this problem.*

*Keywords: reliable computing, numerical filtration, reliability justification.*

## 1. Introduction

We belong to the rather large group of investigators engaged in computational experiments in different fields of science, and as with many of them we have found the problem of reliability justification of results obtained with application of numerical methods very important.

An investigation of scientific problems presents rather difficult processes as far as they apply to the modern technologies of mathematical modelling, numerical experiment and program complex development for this type of experiment. The solving of a multi-parameters problem is accompanied by the appearance of different error types on all these stages. This includes an error of input data, round-off errors, an error of numerical method and an additional unobservable error. That's why the result reliability must be justified. The sources of these errors types are restriction of time, memory capacity, mantissa digits number and reliability.

The analysis of a great number of works on this subject shows that insufficient attention has been paid to the questions of reliability justification of results obtained by mathematical modelling. The verification of numerical data presents great complexity because a considerable part of each investigation (for example a program) is usually not shown. So, the quality of numerical errors analysis depends on researcher experience and intuition.

The comparison of numerical data with physical experiments results often applying to error estimation allows estimation of the approximation error only. But it contains a mathematical model error, an experiment error and an error of computation:

$$\Delta_{\text{approximation}} = \Delta_{\text{model}} + \Delta_{\text{experiment}} + \Delta_{\text{computation}} .$$

Without the computation error estimate this sum doesn't allow estimation of the error of the mathematical model. In this case the sum module can be less than every summand if these summands have different signs. The situation can adversely vary if the problem parameters change. The possibility to estimate the model error appears when independent computation upper error estimate $\overline{\Delta_{\text{computation}}}$ and experiment upper error estimate $\overline{\Delta_{\text{experiment}}}$ are obtained

$$\left|\Delta_{\text{model}}\right| \leq \left|\Delta_{\text{approximation}}\right| + \overline{\Delta}_{\text{experiment}} + \overline{\Delta}_{\text{computation}} .$$

Assuming the model error dependencies have smooth characters, the number of checked physical experiments can be essentially decreased to verify the mathematical model in some range.

Simplified methods of reliability justification of error estimation predominating in computation practice don't possess the required reliability.

Evidently, the earliest method provides a difference of values from distinct experiments as the error estimate:

$$\Delta_z = |z_1 - z_2|. \qquad (1)$$

The values $z_1$ and $z_2$ can be obtained by different numerical methods or by one method with various values of a digitization parameter $n$ (for example, number of grid knots).

If $z_1$ and $z_2$ are approximate values, i.e. $z_1 = z + \Delta_1$, $z_2 = z + \Delta_2$, where $z$ is an accurate value, $\Delta_1$, $\Delta_2$ are unknown errors, then for fulfilling of the inequality:

$$|z_1 - z_2| = |\Delta_1 - \Delta_2| \le \varepsilon,$$

the estimates $|\Delta_1| \le \varepsilon$ and $|\Delta_2| \le \varepsilon$ are valid if and only if the errors $\Delta_1$ and $\Delta_2$ have different signs. In this case, if any sign combination is considered to be equiprobable, then half of the estimates are incorrect. If we take into account the results obtained by one method with variation of the digitization parameter $n$, then the errors most likely have one sign, but the error of value obtained for greater than $n$ can be considerably less than the other error. Then the estimate (1) becomes rather acceptable. Along with the simplicity of this method, it can explain the popularity of such estimations.

So, the estimate (1) can be used only with fulfilment of the defined conditions, which can be violated because of presence of different kinds of error. Therefore, a procedure of the conditions fulfilment verification and estimate justification must exist. But now it is not observed at all. Without such a procedure, errors of estimates of ten to100 times are possible, if the formula (1) is used, which is shown below.

At the beginning of the 20th century the Runge rule of error estimation and the Richardson rule of more precise result definition were worked out on the basis of representation of dependence of an approximate result on $n$ as [1]:

$$z_n = z + c_1 n^{-k} + o\left(n^{-k}\right), \qquad (2)$$

where $z$ is an accurate value; $z_n$ is an approximate result obtained for the mesh points number (or number of sum addends) equal to $n$; $c_1$ is a coefficient independent of $n$; $k$ is an accuracy order of method.

Under the supposition that the small quantity $o\left(n^{-k}\right)$ may be neglected, we have one equation with two unknowns $z$ and $c_1$ for known $k$ (2). For another $n$ ($n_1 = n/Q$) we obtain a similar equation. Solving the system of the two equations the approximate equalities are obtained:

$$z \approx z_n + \frac{z_n - z_{n/Q}}{Q^k - 1}, \quad z_n - z \approx -\frac{z_n - z_{n/Q}}{Q^k - 1}. \qquad (3)$$

The first equality gives the required value $z$ (extrapolated by the Richardson rule). The second one provides an error estimate of the approximate value $z_n$ by the Runge rule.

The mathematical model of an error in more common cases is represented as follows:

$$z_n - z = c_1 n^{-k_1} + c_2 n^{-k_2} + \ldots + c_L n^{-k_L} + o\left(n^{-k_L}\right). \qquad (4)$$

For the finite difference formulas of numerical differentiation and for the Newton–Cotes quadrature formulas of numerical integration the values $k_j$ are integer positive numbers according to the expansion by the Taylor and Euler–Maclaurin formulas [Bakhvalov, Zhidkov & Kobel'kov, 1987].

Let us consider a sequence of the values $n_i = n_1 Q^{i-1}$ and a sequence of numerical results $z_{n_i}, i = 1, \ldots, L+1$, obtained for corresponding different mesh including $n_i$ nodes. $Q = 2$ in most cases. Let us neglect the small quantities and substitute the pairs $n_i, z_{n_i}$ in (4), and then we have system of $L + 1$ linear algebraic equations. This system is the same as the system in the problem of interpolation of some function $f(x)$, $(x = n^{-1})$, by algebraic polynomial $P^L(x)$ given with a table of the values $z_{n_i} = f(x_i)$. We define the required value $z$ by extrapolation of $P^L(x)$ to the value $x = 0$.

If the Aitken recurrent relation [Bakhvalov, Zhidkov & Kobel'kov, 1987] is used to obtain the interpolational polynomial, the problem is reduced to sequential application of the Richardson extrapolation formula (3) to the all pairs of neighbouring values $z_{n_i}$ (for $i$ differing by unit) for $k = k_1$. Then the pairs of the extrapolated values are extrapolated once more for $k = k_2$ and so on.

The Romberg method [Romberg, 1955] consists of the steps mentioned above.

The algorithms $\delta^2$, $\varepsilon$, $\theta$ and others [Aitken, 1926] are applied, if the parameters $k_j$ are unknown in (2), (4). Many other algorithms of the convergence acceleration are worked out for extrapolation and more precise definitions of computation results, in particular for integration of systems of stiff and singular ordinary differential equations and for solving of the Euler and the Navier–Stokes equations and so on.

Note that the possibility of neglecting the last addendum in (2) and (4) is the condition of applicability of these algorithms. It is justified in scientific literature by proof that the last addendum is an infinitesimal quantity $o\!\left(n^{-k_L}\right)$ or $O\!\left(n^{-k_L-1}\right)$. But *n* is always bounded under the conditions of the computer resources restriction. For a finite *n* an infinitesimal quantity $o\!\left(n^{-k_L}\right)$ may not have a really small value in comparison with $c_L n^{-k_L}$. The presence of round-off errors restricts this magnitude below and leads to its growth with increase of *n* in most cases. Also the possibility of error existence in programs isn't to be ignored.

Thus, a strict proof of convergence can't be used as justification of the practical error estimate of numerical results.

Moreover, the estimation methods of partial errors are known (for example, conservation law violation in nonstationary problem solutions). There is a risk of making a mistake in the whole error if the partial estimates are applied.

A lot of published numerical results contain errors in digits declared as valid.

Let us consider the problem of the Stokes solitary wave as an example. This problem has been solved by many investigators [Longuet-Higgins & Fenton, 1974] during recent years. This problem statement is considered in detail at the end of the paper.

The numerical results of the Stokes solitary wave parameters obtained by the different authors are presented in Tables 1 and 2. The errors in units of the last digit (if the error is greater than unit) are given in brackets.

**Table 1.** The values of Froude number (*Fr*) published by different authors.

| Fr | | Reference |
|---|---|---|
| 1.286 | (5) | Longuet-Higgins &Fenton, 1974 [Longuet-Higgins & Fenton, 1974] |
| 1.2909 | | Fox, 1977 [Fox, 1977] |
| 1.290906 | (15) | Hunter & Vanden-Broek, 1983 [Hunter & Vanden-Broek, 1983] |
| 1.290889 | (1) | Williams, 1981 [Williams, 1981] |
| 1.29089053 | (7) | Evans & Ford, 1996 [Evans & Ford, 1996] |
| 1.290890455 | | Sherykhalina, 1995 [Sherykhalina, 1995] |
| 1.2908904558 | | Maklakov, 2002 [Maklakov, 2002] |
| 1.29089045586 | | Zhitnikov & Sherykhalina, 1998 [Zhitnikov & Sherykhalina, 1998] |
| 1.2908904558634 | | Sherykhalina & Zhitnikov, 2001 [Sherykhalina & Zhitnikov, 2001] |
| 1.29089045586335 | | Zhitnikov & Sherykhalina, 2012 [Zhitnikov & Sherykhalina, 2002] |

The authors of the cited papers declare error estimates on the level of unit of the last adduced digits as a rule. The comparison of different authors' results shows that the authors often make mistakes of error estimates five to ten times (at least every second author).

The values of the other parameters of the Stokes solitary wave presented in Table 2 (where *a* is amplitude, *m* – mass, *i* – impulse, *c* – circulation, *t* – kinetic energy, *u* – potential energy; see (20)–(25)) show that all cited authors made mistakes, sometimes in the last two digits. The results of the final column are compared with the results in [Sherykhalina & Zhitnikov, 2001] (Table 3).

**Table 2.** The values of the Stokes solitary wave parameters published by different authors.

| | Longuet-Higgins & Fenton [Longuet-Higgins & Fenton, 1974 ] | Fox [Fox, 1977] | Williams [Williams, 1981] (computed) | Williams [Williams, 1981] (extrapolated) | Evans & Ford [Evans & Ford, 1996] |
|---|---|---|---|---|---|
| *Fr* | 1.286 (−5) | 1.2909 | 1.290891 | 1.290889 (−1) | 1.29089053 (7) |
| *a* | 0.827 (−6) | 0.8332 | 0.833200 | 0.833197 (−2) | 0.833199179 (94) |
| *m* | 1.897 (−73) | 1.968 (−2) | 1.970323 (2) | 1.970319 (−2) | 1.97032019 (−47) |
| *i* | 2.440 (−103) | 2.540 (−3) | 2.543474 (6) | 2.543463 (−5) | 2.54346767 (−47) |
| *c* | 1.653 (−62) | 1.713 (−2) | 1.714569 | 1.714571 (2) | 1.71456873 (−51) |
| *t* | 0.5052(−298) | 0.5339(−11) | 0.535012 (3) | 0.535005 (−4) | 0.535008913 (77) |
| *u* | 0.413 (−25) | 0.4369 (−8) | 0.437675 (2) | 0.437670 (−3) | 0.437672702 (9) |

Note that Williams [Williams, 1981] applied the linear extrapolation of parameters dependence on the residual of a boundary condition for more precise definition of computational results. Table 2 shows that such extrapolation doesn't lead to refinement of the results.

We need to note that the authors of the cited papers applied the generally accepted approach to error estimation and its justification (truly, its absence) and that is why there are no pretensions. And we also do not want to assert that all or most parts of different results published in scientific literature have 1–2 incorrect digits. There are the samples of problem solutions in which a real accuracy corresponds to a declared one. But, as a rule, it is not possible to define where this correspondence takes place and where not, if you have published data only.

It should be noted that the widely accepted mathematical software packages have the same weaknesses. From one side, software developers cannot guarantee reliability of results obtained by a user. From the other side the accompanying documentation has no instructions for justification of obtained results. A software product user has to think of a way of how to do this.

That's why it is possible now to estimate an error in published results correctly only if more precise results appear (or it is known that the errors of two results have different signs). But as Table 1 shows, that chronology of more precise result obtainment can be violated. So, the problem of obtained estimates justification remains relevant for this method of verification too.

With this connection, in this paper a method is proposed allowing at first the determination of a range in which we can confirm the estimates obtained by the Runge rule, the Romberg method and so on. Secondly, estimating with the maximum precision is required: this can be achieved in analysis of results of a concrete numerical experiment and can provide physical reliability of obtained estimates. The physical reliability is guaranteed by determination of an approximate value of a required parameter (called 'standard' below); definition and independent justification of its error estimate (an indeterminacy interval); and checking of intersection fact of intervals obtained by different methods.

## 2. Mathematical Model of Computing Process

The mathematical model of computing process by many numerical methods of some value $z$ can be presented as the function:

$$z_n = z + c_1 f_1(n) + c_2 f_2(n) + \ldots + c_L f_L(n) + \Delta(n), \quad (5)$$

where $z$ is a desired value; $z_n$ is its approximate result obtained for the mesh points number equal to $n$; and $f_1, \ldots, f_L$ are some functions of the mesh points number. All the constants and functions can have both real and complex values.

For the finite difference formulas of numerical differentiation and the Newton–Cotes quadrature formulas $f_j(n) = n^{-k_j}$, where $k_j$ are real numbers as mentioned above.

For the finite difference methods of solving of mathematical physics equations the existence of multiple error components depending on $n$ is confirmed by the computing experiment in [Zhitnikov & Sherykhalina, 2012].

Some numerical methods correspond to the function $f_j(n) = \lambda_j^n$, where $|\lambda_j| < 1$, in particular considered in [Zhitnikov & Sherykhalina, 2013].

The quantity $\Delta(n)$ contains the components that are not included in the sum, the remainder, round-off error and many other components due to both a numerical method and a specific program implementation. Therefore $\Delta(n)$ does not tend to zero with increasing $n$ but becomes greater in most cases.

Let's consider a finite sequence $z_{n_i}^{(0)} = z_{n_i}$, $i = 1, \ldots, I$ of numerical results obtained for grids with $n_i$ nodes. Then we can write the equalities system as such:

$$z_{n_1} = z + c_1 f_1(n_1) + c_2 f_2(n_1) + \ldots + c_L f_L(n_1) + \Delta(n_1),$$

$$z_{n_2} = z + c_1 f_1(n_2) + c_2 f_2(n_2) + \ldots + c_L f_L(n_2) + \Delta(n_2),$$

$$\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots$$

$$z_{n_I} = z + c_1 f_1(n_I) + c_2 f_2(n_I) + \ldots + c_L f_L(n_I) + \Delta(n_I). \quad (6)$$

If $\Delta(n_i)$ is the unknown required parameter along with $z$, $c_1, \ldots c_L$, then the number of the unknowns in the system (6) is always greater than the equations number, and this system of equations has an infinite number of solutions, including an exact one.

We can get an approximate value $\hat{z}$ of the exact solution $z$ and the error estimate applying the regularization

methods [Morozov, 1987]. However, the known methods of regularization require some additional *a priori* information about the unknowns $\Delta(n_i)$, and perhaps about the required parameters $z$, $c_j$. Estimates obtained by such methods are dependent on this additional information, so they can hardly be regarded as error estimates.

In order to avoid incorrectness we propose dividing the problem into two parts: the problem of the mathematical model identification based on the results of the numerical experiment and the testing problem by comparison with known partial solutions or other methods.

The first problem is not to determine the theoretical values $z$, $c_j$, but to decompose $z_n$ into components of the form (5). The mathematical notation of this problem has the same form (6), but $\Delta(n)$ and all other members of (5) have an entirely different meaning, because it is known that $\Delta(n)$ does not include $f_j(n)$, $j = 1, ..., L$ components and the constant.

This problem can be solved approximately by filtration, i.e. by elimination of error components. It makes possible the obtainment of the standard $\hat{z}$ and the error estimate $\overline{\Delta}$. In this case the resulting uncertainty interval $\left[\hat{z} - \overline{\Delta}, \hat{z} + \overline{\Delta}\right]$ can contain or not contain the exact value of $z$, for example, due to an error in the program.

The second problem is testing, i.e. comparing with a known particular exact solution (checking falling into the resulting interval) or with the approximate solution obtained independently by other numerical methods (checking intersection of the uncertainty intervals). This method of additional information application does not affect the estimates obtained previously by independent methods, but only confirms or disproves them.

A theoretical estimation of reliability (confidence probability) of the joint solution of these two problems is presented in [Zhitnikov & Sherykhalina, 1999].

## 3. Numerical filtration

Numerical filtration [Zhitnikov & Sherykhalina, 2012; Zhitnikov & Sherykhalina, 2013] is the sequential removal (elimination) of error components, and the determination of the filtered sequences $z_{n_i}^{(j)}$, $j = 1, ..., L$ ($j$ is a serial number of filtering). For the equation (6) the filtration is reduced to the linear combination $z_{n_i}^{(j)} = \alpha_j z_{n_i}^{(j-1)} + \beta_j z_{n_{i-1}}^{(j-1)}$, where $\alpha_j$ and $\beta_j$ are determined by solving of two equations systems:

$$\alpha_j + \beta_j = 1 , \quad \alpha_j f_j(n_{i-1}) + \beta_j f_j(n_i) = 0 .$$

Then:

$$z_{n_i}^{(j)} = z_{n_i}^{(j-1)} + \frac{z_{n_i}^{(j-1)} - z_{n_{i-1}}^{(j-1)}}{R_j - 1} , \quad R_j = \frac{f_j(n_{i-1})}{f_j(n_i)} , \quad (7)$$

$$c_m^{(j)} = \left(\alpha_j R_m + \beta_j\right) c_m^{(j-1)} = c_m^{(j-1)} \frac{R_j - R_m}{R_j - 1} , \quad m = j + 1, ..., L.$$

The filtration result is the new sequence $z_{n_i}^{(j)}$ without the components $f_j(n)$. This sequence can be subjected to repeated filtration because of preservation of the form of the components that corresponds to fulfilment of the condition $R_j = \text{const}_j$. If $f_j(n) = \lambda_j^n$ (for real or complex $\lambda_j$), then $n_i - n_{i-1} = \alpha = \text{const}_1$, $R_j = \gamma_j^{-\alpha}$, and each of the exponential components is the geometrical progression. If $f_j(n) = n^{-k_j}$, then $n_i / n_{i-1} = Q = \text{const}_2$, $R_j = Q^{k_j}$, and (7) coincides with the Richardson's extrapolation formula. In this case the algorithm of multiple filtrations represents the Romberg method [Romberg, 1955].

Results of multiple filtrations can be presented as a matrix of computed and filtered values:

| $n$ | $z_n$ | $z_n^{(1)}$ | $z_n^{(2)}$ | $z_n^{(3)}$ | ... |
|-----|-------|-------------|-------------|-------------|-----|
| 10 | $z_{10}$ | – | – | – | |
| 20 | $z_{20}$ | $z_{20}^{(1)}$ | – | – | |
| 40 | $z_{40}$ | $z_{40}^{(1)}$ | $z_{40}^{(2)}$ | – | |
| 80 | $z_{80}$ | $z_{80}^{(1)}$ | $z_{80}^{(2)}$ | $z_{80}^{(3)}$ | |
| ... | | | | | |

$$(8)$$

Numerical values given in this matrix must be analysed for error estimation and reliability justification of these estimates.

## 4. Error Estimation, its Justification and Analysis Results Visualization

The data processing results obtained under computation of the integral $\int_0^{\pi/2} \sin x\, dx \approx h \sum_{j=0}^{n-1} \sin\left(\frac{h}{2} + jh\right)$, $h = \frac{\pi}{2n}$ are represented at logarithmic scale in Fig. 1. The bold curves are results of comparison with an accurate solution. Decimal logarithms of relative errors $\delta = \left| \Delta_{n_i}^{(j)} / z_{n_i}^{(j)} \right|$ are put on the ordinate axis, i.e. an accuracy expressed in number of accurate decimal digits. The logarithmic scale on the abscise axis is chosen for power components of error dependency $f_j(n) = n^{-k_j}$; for exponential components the scale is linear. For this choice of the scale every component of the dependence (5) is represented as a straight line.

The lines denoted as *0* in Fig. 1 correspond to an accuracy of computed values of the integral, as *1*, *2*... and so on are an accuracy of results of the first, second and so on filtrations.

The tangent of the line angle inclination (which is equal to $k_1$, $k_2$, ...) varies as a result of every filtration, and an up-shift takes place. This demonstrates the accuracy increase. Such lines behaviour confirms the presence and elimination result of concrete components of the dependence.

Comparison and analysis of several lines location allows the determination of both error estimates and fuzzifying of these estimates. The fuzzifying is an error estimate of the error estimate. The ordinates difference of lines for definite *n* is equal to decimal logarithm of relative fuzzifying of the estimate corresponding to the lower line. So the scale unit corresponds to relative fuzzifying equal to 0.1, two unite to 0.01, etc.

Every line in Fig. 1 has two parts. The first one has inclination corresponding to an exponent of the function. The line behaviour in the second part in contrast to the first one has a chaotic character, which is explained by the prevalence of the error component $\Delta(n)$, the module of which can be approximated by the straight line $y = 16.5 - \frac{1}{2} \lg n$.

Let us consider the results of the Runge rule application (Fig. 1b). It is easy to note that the Runge estimate is close to a real error in the range separated by a half of the scale unit from a level of the irregular error $y = 16.5 - \frac{1}{2} \lg n$. The Runge rule is reduced to comparison of an approximate value with the right neighbouring one in (8). Note, that the Richardson extrapolation application is competitive insofar as it is possible only under prevalence of a component in the form of a power function with a definite exponent. In this case the filtration allows a decrease of the error sensitivity (by the factor of ten or more times) and the comparison with the filtered value provides the same result as with the accurate one. The fuzzifying in this case is very small (less than 0.1) and the estimates are rather precise.

Displacement upwards of second parts of the lines in Fig. 2b (with increase of filtrations number) corresponds to the condition of irregular error prevalence. Difference of the values couples can be used as the error estimate in the range of prevalence of irregular error with changing sign if it is used several times. In the Richardson extrapolation formula this difference is divided into the rather high value $Q^{k_j} - 1$. That is why the difference between approximate and extrapolated values is a small quantity and the Runge estimate gives results with overstated accuracy.
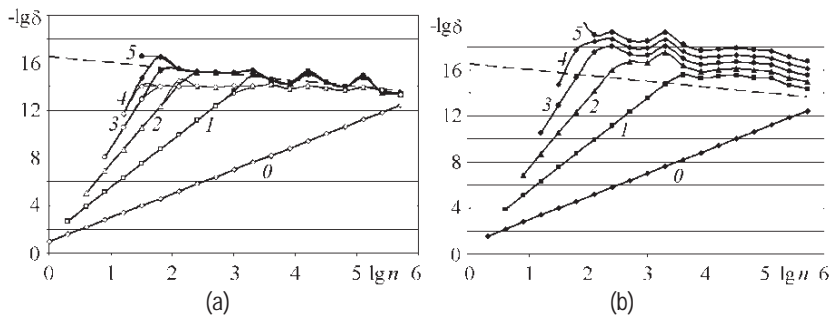


**Fig. 1.** Computation of the integral: (a) comparing with an accurate value and with the standard; (b) the Runge error estimation. The dotted line $y = 16.5 - \frac{1}{2} \lg n$.

Therefore, a comparison of the calculated and the filtered values with a number, called the 'standard', is proposed. The error in the standard choice does not perform the marked disadvantage of the Runge rule ('seeming preciseness'),

caused by the dependence of the estimate on concrete change regularity of error. In the case of standard error existence, the estimate restriction on this error level takes place (Fig. 1b, the thin curves), so the estimates remain correct.

Detection of a locality with approximately constant error allows correction of the standard and increases the result accuracy.

However, the standard choice is inconvenient because of the need for expert intervention in the estimating process online.

## 5. The Standard Selection Rule

To formalize the procedure of the standard selection a two-step method for error estimation is proposed.

The sequence has the form of (6) each time after repeating filtration. Each member of the sequence contains at least two unknowns: the required $z$ and the error.

In order to avoid uncertainty, division of the stages of the error estimation and definition of the required $z$ is suggested. For this the first stage is the filtration by the formula:

$$z_{n_{i-1}}^{(0)} = z_{n_{i-1}} - z_{n_i}, \qquad (9)$$

which excludes the unknown required $z$ from the system (6). Thus, the further filtration by (7) is an error estimation, independent on the standard $z$ choice (Fig. 2a).
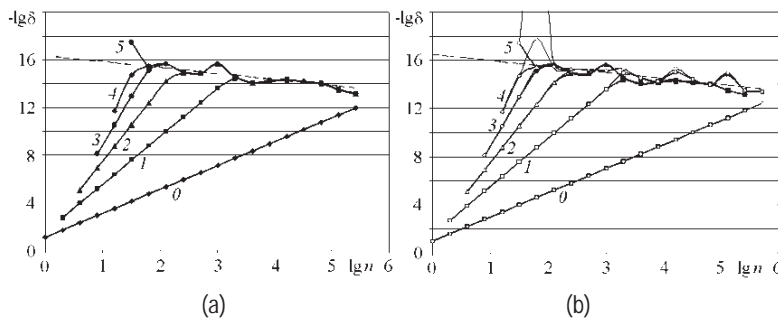


**Fig. 2.** The results of the two-stage error estimate with exclusion of the required value by the formula (9): (a) differences estimate; (b) comparison with the standard (thin lines) and covering the diagram (a) (bold lines).

The estimate obtained by this method allows choice of the best value, from the point of view of an error minimum (or a combination of similar error values), ratio of $n_i$ and $j = j_0$. Thus, the problem is reduced to the minimization problem for some constant $k = 0, 1, 2, ...$ to find minimum $\overline{\Delta}$ for $i$ and $j$ with the constraints:

$$-\overline{\Delta} \le \zeta_{n_i}^{(j)} \le \overline{\Delta},$$

$$\cdots\cdots\cdots\cdots\cdots\cdots \qquad (10)$$

$$-\overline{\Delta} \le \zeta_{n_{i+k}}^{(j)} \le \overline{\Delta}.$$

Values $n_i$ and $j = j_0$ which are obtained by solving the minimization problem are used to determine the values $\hat{z} = z_{n_i}^{(j_0)}$ (the standard) by filtration of the sequence $z_{n_i}^{(0)} = z_{n_i}$ by the formula (7). Superposition of the diagrams for pairwise subtraction (9) and for comparison with the standard allows verification of correctness of the standard choice. Fig. 2 shows that the results coincide with estimated precision.

Note that the transformation (9) changes the components of (5) as such:

$$\zeta_{n_{i-1}}^{(0)} = ... + \left[ f_j(n_{i-1}) - f_j(n_i) \right] c_j^{(j-1)} + ... = ... + \left[ 1 - R_j^{-1} \right] c_j^{(j-1)} f_j(n_{i-1}) + ..., \qquad (11)$$

but for $|R_j| \gg 1$ it is a negligible change (Fig. 2, б).

In the case of considerable distinctions the variation of error components (11) can be taken into account by means

of modification of the transformations (9), (7):

$$\zeta_{n_i}^{(0)} = \frac{z_{n_i} - z_{n_{i+1}}}{1 - R_1^{-1}}, \quad \zeta_{n_i}^{(j)} = \left[\zeta_{n_i}^{(j-1)} + \frac{\zeta_{n_i}^{(j-1)} - \zeta_{n_{i-1}}^{(j-1)}}{R_j - 1}\right]\frac{1 - R_j^{-1}}{1 - R_{j+1}^{-1}} = \frac{z_{n_i}^{(j)} - z_{n_{i+1}}^{(j)}}{1 - R_{j+1}^{-1}}, \quad j = 1, 2, \ldots.$$

Thus the formal procedures of standard determination $\hat{z} = z_{n_i}^{(j_0)}$ and verification of its choice correction are defined. Together with the estimate (10), the interval of uncertainty $\left[\hat{z} - \bar{\Delta}, \hat{z} + \bar{\Delta}\right]$ is determined.

To define the other components of the dependence (1) the identification method offered in [Zhitnikov, Sherykhalina & Porechny, 2010], based on filtration, eliminates the constant from (1) and modifies (1) so that the coefficients $c_1$, $c_2$, etc. appear in the first place in the expansion. Then the sequence determination of these constants is carried out by the two-stage filtration method considered above.

## 6. Formulation and Solution of the Problem of the Stokes Solitary Wave

Let us consider the problem of the Stokes solitary wave, mentioned above, as an example.

The stream of ideal fluid flows along the horizontal rectilinear wall *ADC* (Fig. 3a). The acceleration of gravity points vertically downwards. The solution for a solitary wave of the maximum height is sought. The wave possesses $2\pi/3$ cusp at its crest (the Stokes wave). Asymptotic width of the flow is *h*, and velocity at infinity is $V_\infty$. A pressure *P* on the free surface is equal to the atmospheric pressure $P_0$.

We chose the coordinate system so that the *X*-axis coincides with the straight line *ADC*, and *Y*-axis points vertically upwards and passes through the cusp at the wave crest (the point *B*). We consider the complex variable plane $Z = X + iY$. Let $z = Z/h$.

The free surface shape is unknown. But on the free surface *A'BC'* Bernoulli's equation relates the velocity vector module *V* to the height of point *Y* for $P = P_0$:

$$\left(\frac{V}{V_\infty}\right)^2 + \frac{2}{Fr^2}\frac{Y}{h} = \text{const} = 1 + \frac{2}{Fr^2}, \quad Fr = \frac{V_\infty}{\sqrt{gh}}. \quad (12)$$

Here *g* is gravity acceleration, and *Fr* is the Froude number.
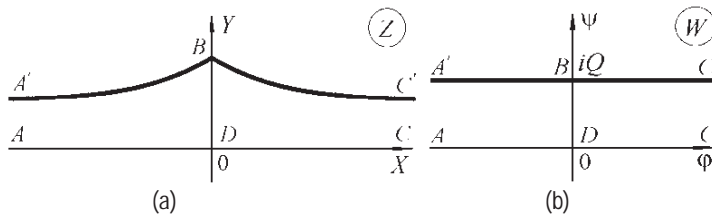


**Fig. 3.** Domain shape corresponding to the flow on the planes: (a) physical; (b) complex potential.

The domain corresponding to the stream on the complex potential plane *W* is a band (Fig. 3b) as far as the boundaries are stream lines.

We choose the band $\chi = \sigma + i\upsilon$ of the width 1 as the parametric variable range (Fig. 4b). In this case the free surface is mapped onto the real axis $\chi = \sigma + i0$.

So, the dependence $W(\chi)$ is determined by the formula:

$$W(\chi) = \varphi + i\psi = Qw(\chi) = Q(-\chi + i), \quad (13)$$

where $\varphi$ is the velocity potential, $\psi$ is a stream function, and $Q = hV_\infty$ is fluid expenditure in the stream.

It is convenient to use the plane of logarithmic hodograph of velocity (the Levi-Chivita function) for problem solving:

$$\omega = i\ln\frac{dw}{dz} = \theta + i\tau. \quad (14)$$

Here $\frac{dw}{dz} = \frac{V}{V_\infty}e^{-i\theta}$ is complex conjugate quantity to dimensionless velocity, and *V* and $\theta$ are module and angle of inclination of velocity vector to *X*-axis, $\tau = \ln(V/V_\infty)$.

*Mediterranean Journal of Social Sciences*
MCSER Publishing, Rome-Italy

The function $\omega(\chi)$ must satisfy the following conditions:

1) the real part $\operatorname{Re}\omega(\chi)=\theta=0$ for Re $\chi=0$, $0<\operatorname{Im}\chi\leq 1$;

2) the real part $\operatorname{Re}\omega(\chi)=0$ for $\operatorname{Im}\chi=1$;

3) the equation (12) connects the real and imaginary parts $\omega(\chi)$ at Im $\chi=0$;

4) the real part $\operatorname{Re}\omega(\chi)\to\pm\dfrac{\pi}{6}$ for $\chi=\sigma+i0$, $\sigma\to\pm0$;

5) the magnitude $\omega(\chi)\to 0$ for $\operatorname{Re}\chi\to\pm\infty$.



**Fig. 4.** Domain shape corresponding to the flow on the planes: (a) logarithmic hodograph of velocity; (b) parametrical variable $\chi$.

According to the symmetry we can consider the right half of the stream.

We represent $\omega$ as a sum:

$$\omega(\chi)=\omega_1(\chi)+\omega_2(\chi)=\theta_1+i\tau_1+\theta_2+i\tau_2,$$

where $\omega_2(\chi)$ is the function satisfying the conditions 1, 2, 4, 5. $\omega_1(\chi)\to 0+i\tau_1(0)$ for $\chi\to 0$ is a result of such representation. So, the function $\omega_1(\chi)$ is continuous on the boundary.

The function $\omega_1(\chi)$ is defined as follows. The solution is sought on the boundary $\chi=\sigma$ at nodes $\sigma_m$ $(m=0,\ldots,n)$. The values $\operatorname{Re}\omega_1(\sigma_m)=\theta_m$ are required. We assume $\operatorname{Re}\omega_1(\sigma_n)=0$ for $\sigma=\sigma_n$, because $\omega_1(\sigma)$ decreases rapidly (as exponent), if $\sigma\to\pm\infty$. Moreover $\theta_0=0$, as far as $\operatorname{Re}\omega_2(0+0)=\pi/6$ (according to the problem conditions). The values $\operatorname{Re}\omega_1(\sigma)$ at points between the nodes are defined with the help of a cubic spline which has two continuous derivatives.

The Schwartz formula [Lavrentiev & Shabat, 1973] is applied for restoration of the function $\omega_1(\chi)$; with regard to the function, $\operatorname{Re}\omega_1(\sigma+i0)$ is an odd function of $\sigma$ and $\operatorname{Re}\omega_1(\sigma+i)=0$:

$$\omega_1(\chi)=-i\int_0^\infty \operatorname{Re}\omega_1(\sigma)\frac{\operatorname{sh}\pi\sigma}{\operatorname{ch}\pi\sigma-\operatorname{ch}\pi\chi}d\sigma.$$

For $\chi\to\sigma_m+i0$ according to the Sokhotskii formula [Lavrentiev & Shabat, 1973]:

$$\omega_1(\sigma_m)=-i\cdot v.p.\int_0^\infty \operatorname{Re}\omega_1(\sigma)\frac{\operatorname{sh}\pi\sigma}{\operatorname{ch}\pi\sigma-\operatorname{ch}\pi\sigma_m}d\sigma+\operatorname{Re}\omega_1(\sigma_m).$$

The calculation of the principal value of the integral for $0\leq m<n$ is practically made by the formula:

$$v.p.\int_0^{\sigma_n}\operatorname{Re}\omega_1(\sigma)\frac{\operatorname{sh}\pi\sigma}{\operatorname{ch}\pi\sigma-\operatorname{ch}\pi\sigma_m}d\sigma=$$

$$\int_0^{\sigma_n}\frac{\operatorname{Re}\omega_1(\sigma)\operatorname{sh}\pi\sigma-\operatorname{Re}\omega_1(\sigma_m)\operatorname{sh}\pi\sigma_m}{\operatorname{ch}\pi\sigma-\operatorname{ch}\pi\sigma_m}d\sigma-\operatorname{Re}\omega_1(\sigma_m)\sigma_m+$$

$$\frac{1}{\pi}\operatorname{Re}\omega_1(\sigma_m)\ln\frac{e^{\pi\sigma_n}-e^{\pi\sigma_m}}{e^{\pi\sigma_n}-e^{-\pi\sigma_m}}. \qquad (15)$$

The two points Gauss quadrature formula is applied for numerical integration between the two knots $\sigma_{m-1}$ and $\sigma_m$. The summand $\omega_2(\zeta)$ is introduced for solution of the singularity isolation at the point $\chi=0$ (Fig. 4a).

Let's present the equation (12) in differential form:

$$e^{3\tau}\frac{d\tau}{d\sigma}-\frac{1}{Fr^2}\sin\theta=0 \qquad (16)$$

The function is determined by:

$$f(\chi)=\frac{1-e^{-\pi\chi}}{\left(1+ie^{-\pi\chi/2}\right)^2}$$

It has real positive values for $\chi=\sigma+i$ and for $\chi=0+i\upsilon$ ($0\le\upsilon\le1$). It tends to 1 for $\mathrm{Re}\chi\to\pm\infty$ and tends to 0 for $\chi\to0$.

The function $\omega_2(\chi)$ is presented in the form:

$$\omega_2(\chi)=\frac{i}{3}\ln f(\chi)+iC_1\left[(f(\chi))^\beta-1\right] \qquad (17)$$

This function satisfies the problem conditions on *AD* and *DB* and it has necessary singularity for $\chi\to0$ as is shown below.

The following estimates take place for $0.5<\beta<1$, $\chi=\sigma+i0$ and $\sigma\to0$:

$$\omega_2(\sigma+i0)=\theta_2(\sigma)+i\tau_2(\sigma)=\frac{\pi}{6}+C_1\sin\frac{\pi\beta}{2}\left(\frac{\pi\sigma}{2}\right)^\beta+\frac{i}{3}\ln\frac{\pi\sigma}{2}-iC_1+$$

$$iC_1\cos\frac{\pi\beta}{2}\left(\frac{\pi\sigma}{2}\right)^\beta+O(\sigma)+O\left(\sigma^{1+\beta}\right)$$

$$\frac{d\tau_2}{d\sigma}=\frac{1}{3\sigma}+C_1\frac{\pi}{2}\beta\cos\frac{\pi\beta}{2}\left(\frac{\pi\sigma}{2}\right)^{\beta-1}+O\left(\sigma^\beta\right)$$

Note, the values $\theta_1(0)=0,\frac{d\tau_1}{d\sigma}(0)=0$, as far as $\theta$ is an odd function and $\tau$ is even function on $\sigma$. The obtained estimates are substituted in (16) and then:

$$\frac{\pi}{6}+C_1\frac{\pi}{2}(\beta+1)\cos\frac{\pi\beta}{2}\left(\frac{\pi\sigma}{2}\right)^\beta+O\left(\sigma^{2\beta}\right)=$$

$$e^{-3(\tau_1(0)-C_1)}\frac{1}{Fr^2}\left[\frac{1}{2}+\frac{\sqrt{3}}{2}C_1\sin\frac{\pi\beta}{2}\left(\frac{\pi\sigma}{2}\right)^\beta+O(\sigma)+O\left(\sigma^{2\beta}\right)\right]$$

Setting equal the values of the same order $O(1)$ and $O(\sigma^\beta)$ we have the expressions:

$$C_1-\tau_1(0)=-\frac{1}{3}\ln\left(\frac{3}{\pi}\frac{1}{Fr^2}\right), \qquad (18)$$

$$(\beta+1)\cot\frac{\pi\beta}{2}=\frac{1}{\sqrt{3}} \qquad (19)$$

The value $0<\beta<1$ is defined from the solution of the equation (19) that was obtained earlier in [Longuet-Higgins & Fox, 1978].

With regard to (13), (14) the free surface shape is determined by numerical integration of the expression:

$$\frac{dZ}{h}=dz=dx+idy=e^{i\omega}dw=-e^{i\omega(\chi)}d\chi$$

The formula of mean rectangles is applied for the numerical integration with the following filtration for achievement of the accuracy up to the order $10^{-16}$.

The problem is solved numerically by the collocation method. The equation (12) is fulfilled at the discrete points of the axis *BA'* ($\sigma_m$, $m=\overline{0,n-1}$). Moreover, the equation (18) must be fulfilled. The system of $n+1$ nonlinear equations is solved numerically with respect to the parameters *Fr*, $C_1$, $\theta_m$ ($m=\overline{1,n-1}$) by the Newton method with step regulation. The sum of residuals squares in the all equations is minimized. The solution search is stopped if the residuals in modulus are less than $10^{-15}$.

The nodes are defined with varying steps according to the formula:

$$\sigma_m=\frac{\sigma_n}{M+\alpha}\left[(M-1)\left(\frac{m}{n}\right)^{\alpha+1}+(\alpha+1)\frac{m}{n}\right]$$

for the values $\sigma_n = 35;40$, $M = 1000$, $\alpha = 3$, chosen experimentally. For $\alpha > 0$ the ratio of nodes density at the end and beginning of the grid is $\dfrac{d\sigma_m/dm(n)}{d\sigma_m/dm(0)} = M$.

## 7. Results analysis

Dimensionless parameters characterizing physical properties of solitary wave [Evans & Ford, 1996] are:

– amplitude
$$a = y_B - 1 = \frac{Fr^2}{2}, \qquad (20)$$

– mass
$$m = \frac{M}{\rho h^2} = \int_{-\infty}^{\infty} [y(x) - 1] dx, \qquad (21)$$

– impulse
$$i = \frac{I}{\rho \sqrt{gh^5}} = \int_{-\infty}^{\infty} dx \int_{1}^{y(x)} \left(1 - \frac{V_x}{V_\infty}\right) dy, \qquad (22)$$

– potential energy
$$u = \frac{U}{\rho g h^3} = \frac{1}{2} \int_{-\infty}^{\infty} (y(x) - 1)^2 dx, \qquad (23)$$

– kinetic energy
$$t = \frac{T}{\rho g h^3} = \frac{1}{2} \int_{-\infty}^{\infty} dx \int_{1}^{y(x)} \left(\left(1 - \frac{V_x}{V_\infty}\right)^2 + \left(\frac{V_y}{V_\infty}\right)^2\right) dy, \qquad (24)$$

– circulation
$$c = \frac{C}{\rho \sqrt{gh^3}} = \int_{-\infty}^{\infty} \left[1 - \sqrt{\left(1 + y'(x)^2\right)\left(1 - \frac{2}{Fr^2}(y(x) - 1)\right)}\right] dx. \qquad (25)$$

In formulas (20)–(25), $\rho$ is fluid density.
According to [Evans & Ford, 1996] the parameters are related as:

$$i = mFr, \quad t = (i - c)\frac{Fr}{2}, \quad u = (Fr^2 - 1)\frac{m}{3}. \qquad (26)$$

We have published the results of this problem in [Sherykhalina & Zhitnikov, 2001] solved by the Levi-Chivita method taking into account the singularities of the solution at the cusp point and at the infinity analogously [Terentiev & Zhitnikov, 2006]. These results are presented in Table 3 (the columns 2, 3). The filtration formula, similar to the Aitken one, applying the least squares method [Sherykhalina & Zhitnikov, 2001; Zhitnikov & Sherykhalina, 2012], is used for more precise definition and error estimation. The obtained estimates allow declaration of accuracy up to the 13th digit for the Froude number and up to two to three units of the 13th digit for the other parameters (Table 3, column 3). Later on [Sherykhalina & Zhitnikov, 2001], the ε-algorithm is applied to the data obtained by the same method. The dependencies $z_{n_i}$ diagrams covering takes place for filtration results analysis for different sequences $n_i$ beginning from 5, 6, 7, 8 and further multiple doubling of these initial values. This allows definition of a more precise value of the Froude number up to 14 digits (the last line of Table 1) and the mass value up to the unit of the 13th digit. This confirms the estimates of these parameters in [Sherykhalina & Zhitnikov, 2001; Zhitnikov & Sherykhalina, 2012]. But it is necessary to use another method for independent verification of these estimates. The method considered in this part is developed for this purpose.

The results calculated by the given method and filtered by the rules described above are presented in Table 3 (the columns 4, 5).

**Table 3.** The values of the solitary wave parameters obtained by different methods.

| | Sherykhalina & Zhitnikov [Sherykhalina & Zhitnikov, 2001] (computed) | Sherykhalina & Zhitnikov[Sherykhalina & Zhitnikov, 2001] (filtered) | The results of the current paper (computed) | The results of the current paper (filtered) |
|---|---|---|---|---|
| *Fr* | 1.290890455863 | 1.2908904558634 | 1.2908904558635 | 1.290890455863341 |
| *a* | 0.833199084519 | 0.8331990845196 | 0.8331990845197 | 0.833199084519532 |
| *m* | 1.9703206602 | 1.9703206601317 (−2) | 1.97032066009 | 1.97032066013185 |
| *i* | 2.5434681352 | 2.5434681351545 | 2.54346813510 | 2.54346813515454 |
| *c* | 1.7145692406 | 1.7145692405337 (2) | 1.7145692405 | 1.71456924053350 |
| *t* | 0.5350088360 | 0.5350088359709 | 0.53500883596 | 0.53500883597099 |
| *u* | 0.43767269344 | 0.4376726934439 | 0.437672693441 | 0.43767269344390 |

Let us analyse and justify the error estimates of the obtained results. Combination of the Aitken and Richardson formulas is used as the filtration formula. The values $k_j$ are defined by the Aitken formula $$Q^{k_j} \approx \frac{z_{n_{i-1}}^{(j-1)} - z_{n_{i-2}}^{(j-1)}}{z_{n_i}^{(j-1)} - z_{n_{i-1}}^{(j-1)}}$$ for $Q=2$. This makes it possible to find the values $k_j = 4, 5, \cdots$. Then the formula (7) is applied for $R_j = 2^{k_j}$. The ε-algorithm is used for independent verification. The estimates obtained by the two methods were compared.

The results of the filtration and error estimation are presented in Figs. 5 and 6. As described above, dependencies of the error estimates logarithms on logarithm of collocations points number $n$ are represented.
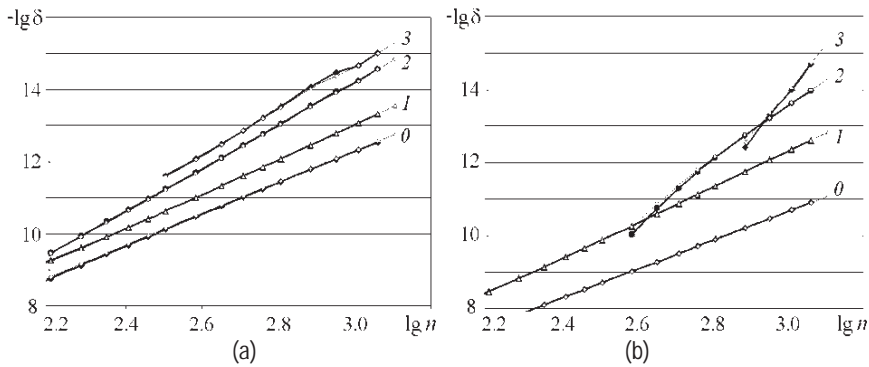


**Fig. 5.** Solution filtration results of the Stokes solitary wave problem: (a) the Froude number; (b) the potential energy.
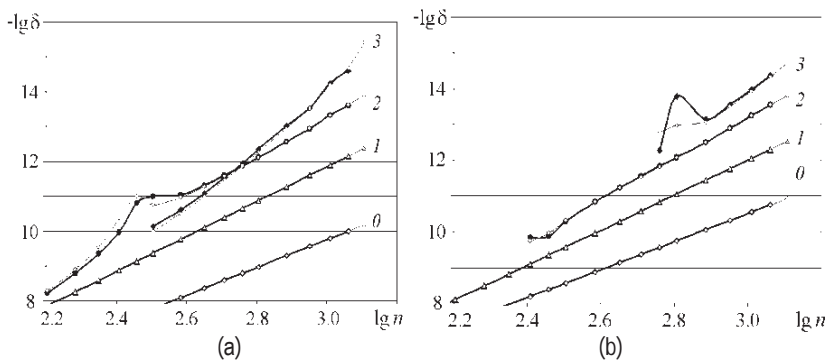


**Fig. 6.** Solution filtration results of the Stokes solitary wave problem: (a) circulation; (b) mass.

The diagrams covering the dependencies $z_{n_i}$ obtained for different sequences $n_i$ beginning from 5, 6, 7, 8, 9 and further multiple doubling of these initial values (up to $n = 1280$) is used for more detailed investigation of the dependencies. The filtration of the every sequence is carried out independently, but the values with more closed $n_i$ are used for pairwise subtraction (9). Such improvement of the filtration method allows detection of a small deviation of lines behaviour corresponding to the pairwise subtraction and comparison with the standard, which requires additional analysis. The upper lines (number 3) are used for an approximate estimate of the level of irregular error $\Delta(n)$ and error estimate fuzzifying but not for error estimates of obtained data.

According to these estimates, the results given in Table 3 (column 5) have the confirmed accuracy of up to two to three units of the 15th digit for the Froude number and the amplitude and up to two to three units of the 14th digit for the other parameters.

## 8. Conclusion

Thus, three problems are useful for error estimation from the practical point of view. The first one is the definition of the range of digitization parameter $n$ and the error $\delta$, where the error estimation by more simple rules (as Runge and Romberg estimations) gives acceptable results. The second problem is to provide the highest possible accuracy of the standard for a given experiment and its error estimate. The third problem is investigation of error components for search mistakes in the program and improvement of numerical algorithms. All these problems are solved by the filtration methods.

A procedure including two stages has been suggested: the identification of a mathematical model, by means of multicomponent analysis of numerical experiment results, and testing.

The filtration realized at the stage of postprocessor treatment of numerical experiment data allows the obtainment of error estimates and essentially raises the efficiency of numerical algorithms. Testing with help of accurate or approximate (obtained by the other method) particular solutions makes possible the confirmation or disproval of these estimates. So the reliability of numerical data and conclusions based on the estimates increases.

The investigation of error dependencies on digitization number $n$ (for example, collocation points number) allows formulation of the conclusion of three areas. At first one can see the area of chaos caused by rough digitization; then the area of prevalence of regular error components, which can be filtered. And then there is the area of the round-off error predominance. Boundaries of the areas are determined experimentally.

The solving of the Stokes solitary wave problem by the new method allows the confirmation of the estimates obtained earlier in [Sherykhalina & Zhitnikov, 2001] with help of the filtration method.

## 9. Acknowledgments

## References

1. Bakhvalov, N. S., Zhidkov, N. P., Kobel'kov, G.M. 1987 *Numerical methods.* Moscow: Nauka.
2. Romberg, W. (1955) Vereinfachte numerische Integration. *Det. Kong. Norkse Videnskabernes Selskabs Forhandlinger,* 28(7), 30-36.
3. Aitken, A. C. (1926) On Bernoulli's numerical solution of algebraic equations. *Proc. R. Soc. Edinb,* 46, 289-305.
4. Wynn, P. (1956) On a device for computing the em(Sn) transformation. *MTAC* 10, 91-96. [Online] Available: http://dx.doi.org/10.2307/2002183
5. Smith, D. A. & Ford, W. F. (1979) Acceleration of linear and logarithmic convergence. SIAM J. Numer. Anal. 16, 223-240. (DOI: 10.1137/0716017.)
6. Smith, D. A. & Ford, W. F. (1982) Numerical comparisons of non-linear convergence accelerations. *Math. Compu,* 38, 481-499. (DOI: 10.1090/S0025-5718-1982-0645665-1.)
7. Longuet-Higgins, M. S. & Fenton, J. D. (1974) On the mass, momentum, energy and circulation of the solitary wave. *Proc. R. Soc.,* 340, 471-493. (DOI: 10.1098/rspa.1974.0166.)
8. Fox, M. J. H. (1977) *Nonlinear effects in surface gravity waves on water.* PhD thesis. Cambridge Univ.
9. Hunter, J.K. & Vanden-Broek, J.-M. (1983) Accurate computations for steep solitary waves. J. Fluid. Mech., 136, 63-71. DOI: (10.1017/S0022112083002050.)
10. Williams, J. M. (1981) Limiting gravity waves in water of finite depth. *Phil.* Trans. R. Soc. Lond., A 302, 139-188. (DOI: 10.1098/rsta.1981.0159.)
11. Evans, W. A. B. & Ford, M. J. (1996) An exact integral equation for solitary waves (with new numerical results for some internal properties). *Proc. R. Soc. Lond., A* 452, 373-390. (DOI: 10.1098/rspa.1996.0020.)
12. Sherykhalina, N. M. (1995) *Development of computational algorithms of hydrodynamical problems solution.* VINITY 2550-B95.
13. Maklakov, D. V. (2002) Almost-highest gravity waves on water of finite depth. *Euro J. Appl. Math.,* 13, 67-93. (DOI: 10.1017/S0956792501004739.)
14. Zhitnikov, V. P. & Sherykhalina, N. M. (1998) The solitary waves shape calculation by numerical-analytical methods. *Fizika volnovyh processov & radiotechnicheskie sistemy,* 1, 103-107.
15. Sherykhalina, N. M. & Zhitnikov, V. P. (2001) Application of extrapolation methods of numerical results for improvement of hydrodynamics problem solution. *Comput. Fluid Dyn. J.,* 10, 309-314.
16. Zhitnikov, V. P. & Sherykhalina, N. M. (2012) *Multicomponent analysis of numerical results.* Saarbrucken: Lambert Academic Publishing.

17. Zhitnikov, V. P. & Sherykhalina, N. M. (2013) Multi-stage filtration of numerical problems solution by methods of complex variable functions theory. *Comput. Tech.,* 18, 15-24.
18. Morozov, V. A. (1987) *Regular methods of ill-posed problems solving* (in Russian). Moscow: Nauka.
19. Zhitnikov, V. P. & Sherykhalina, N. M. (1999) Certainty estimation of numerical results in the presence of several methods of problem solution. *Comput. Tech.,* 4, 77-87.
20. Zhitnikov, V. P., Sherykhalina, N. M. & Porechny, S.S. (2010) Identification problem solving applied to the numerical results estimation. *Sci. Tech. Sheets SPbTU,* 1, 60-63.
21. Lavrentiev, M. A. & Shabat, B. V. (1973) *Methods of the theory of complex variable functions*. Moscow: Nauka.
22. Longuet-Higgins, M. S. & Fox, M. J. H. (1978) Theory of the almost-highest wave. Part 2. Matching and analytical extension. *J. Fluid Mech.,* 85, 769-786. (DOI: 10.1017/S0022112078000920.)
23. Terentiev, A. G. & Zhitnikov, V. P. (2006) Stationary two-dimensional inviscid flow with flexible boundaries including the effect of surface tension. *J. Eng. Math,.* 55, 111-126. (DOI: 10.1007/s10665-005-9021-2.)