# A strongly $S$-stable method
# for solving stiff systems
# of ordinary differential equations

S.A. Gusev

A strongly $S$-stable (by A. Prothero and A. Robinson) one-step noniterated method is presented. Results of numerical calculation showing the advantage of the proposed method in comparison with a similar $L$-stable method are given.

## 1. Introduction

It is known that in the using of $A$-stable methods for the solution to very stiff systems of ordinary differential equations (ODE)

$$\dot{y} = f(y, t), \quad y(0) = y_0, \quad t \geq 0 \tag{1}$$

the stable numerical solutions are not always obtained. Sometimes the accuracy of an obtained numerical solution is less than that is supposed from consistency equations. It make A. Prothero and A. Robinson to enter other concepts of stability of numerical methods such as $S$-*stability* and *strong $S$-stability* [1]. In this connection they considered the following test equation:

$$\dot{y} = \dot{g}(t) + \lambda(y - g(t)), \quad y(0) = y_0, \quad t \geq 0, \tag{2}$$

where $g$ is any defined sufficiently differentiable function, $\lambda$ – complex constant with $\operatorname{Re}(\lambda) \ll 0$.

According to [1] the one-step numerical method is called $S$-*stable*, if in applying it to equation (2) and for any positive constant $\lambda_0$, there exists $h_0 > 0$ such that

$$\left| \frac{y_{n+1} - g(t_{n+1})}{y_n - g(t_n)} \right| < 1, \tag{3}$$

provided $y_n \neq g(t_n)$, for all $0 < h < h_0$ and all $\lambda \in C$ such that $\operatorname{Re}(-\lambda) \geq \lambda_0$, and $t_n, t_{n+1} \in [0, T]$ (any $0 < T < \infty$).

An $S$-stable one-step method refers to as *strongly $S$-stable*, if

$$\frac{y_{n+1} - g(t_{n+1})}{y_n - g(t_n)} \to 0, \tag{4}$$

under $\operatorname{Re}(-\lambda) \to \infty$ for all positive $h$ such that $t_n, t_{n+1} \in [0, T]$.

The advantages of $S$-stable methods in comparison with methods, which have not this property, were evidently demonstrated in the papers [1, 2].

Under the given initial condition $y(0) = y_0$ the function

$$y(t) = e^{\lambda t}(y(0) - g(0)) + g(t) \qquad (5)$$

is the solution to equation (2).

In [1] the equation for the error $e_n = y_n - g(t_n)$ was also considered

$$e_{n+1} = \alpha(z)e_n + h\beta(h, z, g), \qquad (6)$$

where $z = (h\lambda)^{-1}$.

For one-step methods $\alpha(z)$ coincides with the stability function that is obtained as a result of application of these methods to the scalar test equation

$$\dot{y} = \lambda y, \quad y(0) = y_0, \quad t \geq 0. \qquad (7)$$

As it is shown in [1] for $S$-stablility of an one-step method it is necessary and sufficiently that it is $A$-stable and the function $\beta(z)/(1 - |\alpha(z)|)$ is bounded for all $z$ with $0 < \text{Re}(-z) < \tilde{z}$ (any $\tilde{z} > 0$) and for all $g$ with $\dot{g}$ defined and bounded in $[t_n, t_{n+1}]$.

An $S$-stable method is strongly $S$-stable if and only if [1]

$$\lim_{\text{Re}(z)<0, |z|\to 0} \alpha(z) = 0, \qquad (8)$$

$$\lim_{\text{Re}(z)<0, |z|\to 0} \beta(z) = 0. \qquad (9)$$

A method is called as *stiffly accurate* if condition (9) is realised.

Strongly $S$-stable methods are the implicit Euler method and the Rosenbrock type method of the first order

$$y_{n+1} = y_n + [I - hf_y(y_n t_n)]^{-1} f(y_n, t_{n+1}). \qquad (10)$$

As is shown in [1] strongly $S$-stable implicit methods are also some Runge-Kutta methods based on the Radau and Lobatto quadrature formulas.

In this paper a strongly $S$-stable one-step noniterated method with two calculations of right-hand side and one calculation of the Jacobi matrix at one integration step is offered. This method belongs to so called class of $(m, k)$ methods proposed in [3] and it is a four stage $(m, k)$ method of the kind

$$y_{n+1} = y_n + \sum_{i=1}^{4} p_i k_i, \qquad (11)$$

$$[I - ahf_y]k_1 = f(y_n, t_n + \gamma_1 h),$$

$$[I - ahf_y]k_2 = k_1,$$

$$[I - ahf_y]k_3 = f(y_n + \beta_{31}k_1 + \beta_{32}k_2, t_n + \gamma_3 h),$$

$$[I - ahf_y]k_4 = k_3 + \alpha_{42}k_2,$$

where $I$ – unit matrix, $f_y$ – the Jacobi matrix at the point $(y_n, t_n)$, $a$, $p_1$, $p_2$, $p_3$, $p_4$, $\gamma_1$, $\gamma_3$, $\beta_{31}$, $\beta_{32}$, $\alpha_{42}$ are numerical parameters.

## 2. Parameters and properties of the method

For determination of parameters of a method (11) we shall consider consistency conditions of the third order and strong $S$-stability. The third order concistency equations are

$$p_1 + p_2 + p_3 + p_4 + p_5 = 1, \tag{12}$$

$$p_1 a + 2p_2 a + p_3(a + \beta_{31} + \beta_{32}) + p_4(2a + \beta_{31} + \beta_{32}) + 3p_5 a = \frac{1}{2}, \tag{13}$$

$$(p_1 + p_2 + p_5)\gamma_1 + (p_3 + p_4)\gamma_3 = \frac{1}{2}, \tag{14}$$

$$(p_1 + p_2 + p_5)\gamma_1^2 + (p_3 + p_4)\gamma_3^2 = \frac{1}{3}, \tag{15}$$

$$6p_1 a\gamma_1 + 12p_2 a\gamma_1 + 18p_5 a\gamma_1 + 6p_3(a\gamma_3 + (\beta_{31} + \beta_{32})\gamma_1) + $$
$$6(p_4(2a\gamma 3 + (\beta_{31} + \beta_{32})\gamma_1) = 1, \tag{16}$$

$$6p_1 a^2 + 18p_2 a^2 + 36p_5 a^2 + 6p_3(a^2 + 2\beta_{31}a + 3\beta_{32}a) + $$
$$6p_4(3a^2 + 3\beta_{31}a + 4\beta_{32}a) = 1, \tag{17}$$

$$(p_3 + p_4)(\beta_{31} + \beta_{32})^2 = \frac{1}{3}, \tag{18}$$

$$\gamma_3(p_3 + p_4)(\beta_{31} + \beta_{32}) = \frac{1}{3}. \tag{19}$$

In equations (12)–(19) for convenience we denote $p_5 = p_4\alpha_{42}$.

Next we consider equations for parameters providing strong $S$- stability. By application of (11) to equation (2) we have

$$y_{n+1} = y_n + \frac{1}{z}\bar{p}^T B^{-1}(y_n - g(t_n))\bar{e} + h\bar{p}^T B^{-1}\bar{g}' + \frac{1}{z}\bar{p}^T B^{-1}(g(t_n)\bar{e} - \bar{g}), \tag{20}$$

$$\bar{p} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix}, \quad \bar{e} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad \bar{g}' = \begin{bmatrix} \dot{g}(t_n + \gamma_1 h) \\ \dot{g}(t_n + \gamma_1 h) \\ g(t_n + \gamma_3 h) \\ g(t_n + \gamma_3 h) \end{bmatrix}, \quad \bar{g} = \begin{bmatrix} g(t_n + \gamma_1 h) \\ g(t_n + \gamma_1 h) \\ g(t_n + \gamma_3 h) \\ g(t_n + \gamma_3 h) \end{bmatrix},$$

$$B = \begin{bmatrix} \frac{z-a}{z} & 0 & 0 & 0 \\ 0 & \left(\frac{z-a}{z}\right)^2 & 0 & 0 \\ -\frac{\beta_{31}}{z} & -\frac{\beta_{32}}{z} & \frac{z-a}{z} & 0 \\ -\frac{\beta_{31}+\alpha_{42}}{z} & -\frac{\beta_{32}}{z} & 0 & \left(\frac{z-a}{z}\right)^2 \end{bmatrix}.$$

Using equality (20) we obtain the expressions $\alpha(z)$, $\beta(z, h)$, specifying a step error change (6) for our method

$$\alpha(z) = 1 + \frac{1}{z}\bar{p}^T B^{-1}\bar{e}, \tag{21}$$

$$\beta(z,h) = \frac{1}{h}(g(t_n) - g(t_{n+1}) + \frac{1}{z}\bar{p}^T B^{-1}(g(t_n)\bar{e} - \bar{g}) + h\bar{p}^T B^{-1}\bar{g}'). \tag{22}$$

Limiting equalities (8), (9) for (21), (22) give equations on parameters ensuring strong $S$-stability.

By equating to zero in (21) the limit $\alpha(z)$ at $|z| \to 0$, we obtain the first equation that is also a condition of $L$-stability

$$a^2 + p_1 a + p_3(a - \beta_{31}) = 0. \tag{23}$$

Similarly, by equating to zero the limit $\beta(z,h)$ in (22) at $|z| \to 0$, we obtain equality, ensuring the stiff accuracy

$$g(t_n + h) - g(t_n) - \frac{p_1}{a}(g(t_n + \gamma_1 h) - g(t_n)) - \frac{p_3}{a}(g(t_n + \gamma_3 h) - g(t_n)) +$$

$$\frac{p_3 \beta_{31}}{a^2}(g(t_n + \gamma_1 h) - g(t_n)) = 0. \tag{24}$$

Equality (24) is possible, if in Taylor-series expansions $g(t_n + h)$, $g(t_n + \gamma_i h)$, $i = 1,3$ a sum of factors at each degree of $h$ is equal to zero. Thus we have the infinite system of equations

$$1 - \frac{p_1\gamma_1^k}{a} - \frac{p_3}{a}\left(\gamma_3^k - \frac{\beta_{31}\gamma_1^k}{a}\right) = 0, \quad k = 1, 2, \dots . \tag{25}$$

So that the method (11) had the third order and was strongly $S$-stable, the fulfilment of conditions (12)–(19), (23), (25) is necessary. It is possible to show that this system of equations is inconsistent. Indeed, it follows from equations (18), (19) that

$$\beta_{31} + \beta_{32} = \gamma_3, \tag{26}$$

$$p_3 + p_4 = \frac{1}{3\gamma_3^2}. \tag{27}$$

From (27) and (15) we have $(p_1 + p_2 + p_5)\gamma_1^2 = 0$.

Assume $p_1 + p_2 + p_5 = 0$, then from (12) it follows $p_3 + p_4 = 1$. Having substituted the latter equality in (18), (19) and having excluded $\beta_{31} + \beta_{32}$ we obtain $\gamma_3 = 1/\sqrt{3}$. But in division (15) on (14) it is obtained $\gamma_3 = 2/3$.

If assume $\gamma_1 = 0$, then from equalities (14), (15) it follows that $\gamma_3 = 2/3$. But at this condition the stiff accuracy (25) can be executed, if $\gamma_3 = 1$.

So, the system (12)–(19), (23), (25) is inconsistent, but it appears the joint system of equations to be consistent that differs from this system by only in the right-hand side of equation (19), corresponding to differential $f_{ty}f$, precisely, instead of $\frac{1}{3}$ is $\frac{1}{6}$

$$\gamma_3(p_3 + p_4)(\beta_{31} + \beta_{32}) = \frac{1}{6}. \tag{19'}$$

By solution of system (12)–(18), (19'), (23), (25) the following significances of parameters were obtained:

$$a = \frac{1}{3}, \quad p_1 = \frac{1}{3}, \quad p_2 = \frac{19}{12}, \quad p_3 = 0, \quad p_4 = \frac{3}{4},$$
$$\gamma_1 = 1, \quad \gamma_3 = \frac{1}{3}, \quad \beta_{31} = \frac{22}{27}, \quad \beta_{32} = -\frac{4}{27}, \quad \alpha_{42} = -\frac{20}{9}. \tag{28}$$

Thus, the strongly $S$-stable numerical method "almost" of third order is constructed. For nonautonomous ODE systems having zero differential of third order $f_{ty}f$ and for autonomous ODE systems the order of the method is equal to three.

In [1] the concept of *stiff order* of an one-step method was introduced. The stiff order is a pair of integers $(s, r)$, determining orders of asymptotic behaviour of local truncation error $l_n$ at $h \to 0$, $\text{Re}(-h\lambda) \to \infty$, when the method is applied to equation (2). The one-step method has the stiff order $(s, r)$, if at at $h \to 0$, $\text{Re}(-h\lambda) \to \infty$

$$l_n \propto h^{s+1}\lambda^r. \tag{29}$$

Next we determine the stiff order of the method (11) with parameters (28). Assuming in (6) $e_n = 0$ we obtain its local truncation error

$$l_n = h\beta(z, h) = g(t_n) - g(t_{n+1}) + \frac{1}{z}\bar{p}^T B^{-1}(g(t_n)\bar{e} - \bar{g}) + h\bar{p}^T B^{-1}\bar{g}'. \tag{30}$$

From (30) expanding in Taylor-series $g(t_{n+1})$, $\bar{g}$, $\bar{g}'$ it is possible to obtain representation $l_n$ as series on $h$

$$l_n = \sum_{i=1}^{\infty} c_i h^i, \tag{31}$$

with

$$c_1 = g'(t_n)(-1 + \bar{p}^T B^{-1} z^{-1}(z\bar{e} - \bar{\gamma})), \tag{32}$$

where $\bar{\gamma} = [\gamma_1, \gamma_1, \gamma_3, \gamma_3]^T$.

Substituting in (32) significances of factors (28) we obtain

$$c_1 = \frac{2zg'(t_n)}{81z^4 - 108z^3 + 54z^2 - 12z + 1}. \tag{33}$$

Under $|z| \to 0$, $c_1 = O(z)$, and the stiff order of the proposed method is equal $(0, -1)$.
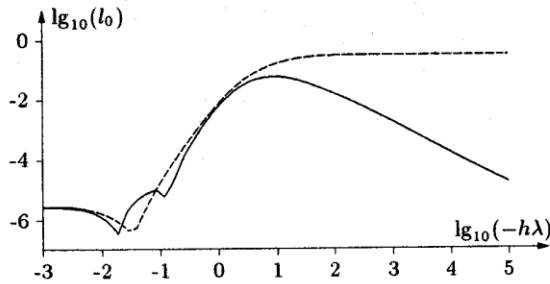
*S.A. Gusev*

## 3. Numerical results

The aim of this section is to demonstrate the advantage of a strongly $S$-stable method in numerical solution of stiff nonlinear ODE systems. For comparison we consider the $L$-stable method of the third order (11) which is not strongly $S$-stable. Following parameters of this method satisfying equations (12)–(19) and (23) were obtained

$$a = \frac{1}{2}, \quad p_1 = \frac{3}{2}, \quad p_2 = -\frac{7}{4}, \quad p_3 = 1, \quad p_4 = -\frac{1}{4},$$

$$\gamma_1 = 0, \quad \gamma_3 = \frac{2}{3}, \quad \beta_{31} = 1, \quad \beta_{32} = -\frac{1}{3}, \quad \alpha_{42} = -2. \tag{34}$$

Let denote $SST$ the method (11) with parameters (28) and $LST$ the method (11) with parameters (34). The local truncation errors $l_0$ of methods $SST$ and $LST$, when they apply to solving with $h = 0.1$ of the test problem from [1]

$$\dot{y} = \dot{g}(t) + \lambda(y - g(t)), \qquad y(0) = 0, \tag{35}$$
$$g(t) = 10 - (10 + t)e^{-t}, \qquad \lambda \in R,$$

are demonstrated in Figure 1. One can see that the local truncation error of the strongly $S$-stable method tends to zero for large $-\lambda h$, that is not so for the $L$-stable method.



**Figure 1.** Local truncation errors $l_0$, when problem (35) solved by $SST$ and $LST$ methods: —— $SST$, — — $LST$

The comparison of these two methods was also made on the solution of an ODE system, which is obtained by application of the Galerkin method with piecewise linear basic functions to parabolic one-dimensional initial boundary value problem with known exact solution [4]

$$C_0 \Psi(T) \frac{\partial T}{\partial t} = K_0 \frac{\partial}{\partial x}(\Psi(T)) \frac{\partial T}{\partial x} - a\Phi(T), \quad t > 0, \ 0 < x < L,$$

$$t = 0 : T = T_0, \tag{36}$$

$$x = 0 : K_0 \Psi(T) \frac{\partial T}{\partial x} = Q(t), \qquad x = L : K_0 \Psi(T) \frac{\partial T}{\partial x} = 0.$$
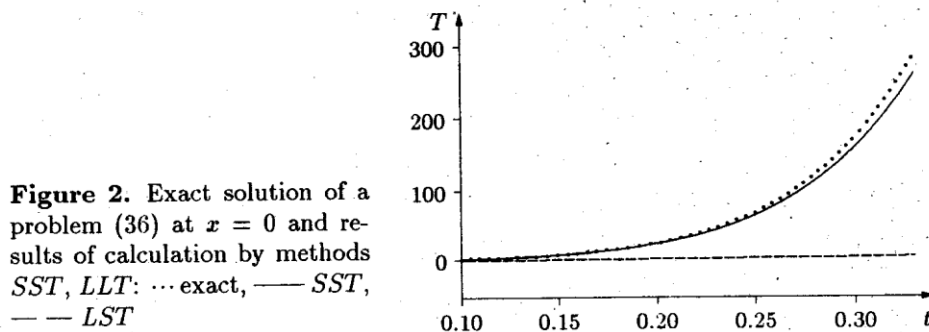
In problem (36) $Q(t)$ is given in the form

$$Q(t) = \frac{Q_2 t}{t_2^2} e^{(-t/t_2)} - \frac{Q_1 t}{t_1^2} e^{(-t/t_1)}, \quad Q_1 > Q_2 > 0, \quad 0 < t_1 < t_2. \quad (37)$$

The functions $\Psi$, $\Phi$ satisfy the condition

$$0 \le \Psi(\theta) = \Phi'(\theta), \quad \Phi(1) = 1.$$

In this case $\Phi(\theta) = \ln(\theta) + 1$, $\Psi(\theta) = 1/\theta$.



**Figure 2.** Exact solution of a problem (36) at $x = 0$ and results of calculation by methods $SST$, $LLT$: $\cdots$ exact, —— $SST$, — — $LST$

The exact solution to problem (36), (37) is represented as series, which complete expression is given in [4]. In calculations the infinite sum of series is replaced by a final sum with a relative error $10^{-6}$. In problem (36) the following significances of constants $C_0 = 1$, $K_0 = 1$, $L = 1$, $Q_1 = 100$, $Q_2 = 1$, $t_1 = 1$, $t_2 = 8$, $T_0 = 1$ were taken. The exact solution to problem (36) at the point $x = 0$ (Exact) and corresponding numerical solutions obtained by the $SST$ and $LST$ methods with the constant step $h = 10^{-3}$ are shown in Figure 2.

# 4. Conclusion

By using the test equation (2) the one-step strongly $S$-stable numerical method is constructed. Numerical results show, the obtained method exeeds in stability properties close to it in a structure $L$-stable method. It is confirmation of the results of the papers [1], [2] about advantage of $S$-stable and strongly $S$-stable methods in solution of stiff nonlinear problems in comparison with methods which have not these stability properties.

# References

[1] A. Prothero, A. Robinson, *On the stability and accuracy of one-step methods for solving stiff systems of ordinary differential equations*, Math. Comp., **28**, 1974, 145–162.

[2] J.G. Verver, *S-stability properties for generalized Runge-Kutta methods*, Numer. Math., 27, 1977, 359–370.

[3] E.A. Novikov, Yu.A. Shitov, *Some methods for solving stiff systems induced by one and two calculations of right hand side*, Mathematical models and solution methods for problems of continuous medium mechanics, Krasnoiarsk, 1986, 11–18.

[4] S.A. Gusev, O.A. Makhotkin, *The decision of problems radiative conductive transfer by a Monte Carlo method*, IV. Approximation of temperature by B-splines in one-dimensional heat conductivity equation, Theory and application of statistical modeling, Novosibirsk, 1991, 70–84.