# Mathematica Balkanica

# Economic Modification of a Sweep's Method for Solution of Tridiagonal Linear Algebraic System of Equation

*J. A. Zarnan*

*Presented by P. Kenderov*

In the present paper we continue the study of P.V. J a n j g a v a from [1] on Sweep's method for solution on tridiagonal systems.

We shall start with the scalar case.

## 1. Scalar case

Let it be given the following special linear tridiagonal system

(1)
$$
\begin{aligned}
bx_1 + cx_2 \quad\quad &= d_1 \\
ax_1 + bx_2 + cx_3 &= d_2 \\
\cdots\cdots\cdots\cdots \quad\quad &\cdots\cdots\cdots\cdots \\
\cdots\cdots\cdots\cdots \quad\quad &\cdots\cdots\cdots\cdots \\
ax_{n-1} + bx_n &= d_n
\end{aligned}
$$

As known [2] the Sweep's method for solution of the system (1) can be applied with two steps: forward and backward. The forward step is to transformate the first $(n-1)$ equations of (1) in equations of the form

$$x_1 = \alpha_1 x_2 + \beta_1$$

$$x_2 = \alpha_2 x_3 + \beta_2$$

$$\cdots\cdots\cdots\cdots\cdots$$

$$\cdots\cdots\cdots\cdots\cdots$$

$$x_{n-1} = \alpha_{n-1} x_n + \beta_{n-1}$$

where the coefficients $\alpha_i$ and $\beta_i$ are obtained by the formulas

$$\alpha_i = -\frac{c}{\alpha_{i-1}a + b}; \quad \beta_i = \frac{d_i - a\beta_{i-1}}{\alpha_{i-1}a + b},$$

where $i = 1, 2, ..., n - 1$, $\alpha_0 = \beta_0 = 0$.

The backward step is to find $x_n$, $x_{n-1}$,...,$x_1$ by the formulas

$$x_n = \frac{d_n - a\beta_{n-1}}{\alpha_{n-1}a + b}$$

$$x_{n-1} = \alpha_{n-1}x_n + \beta_{n-1}$$

$$.....................$$

$$.....................$$

$$x_1 = \alpha_1 x_2 + \beta_1$$

Now let the matrix of the system (1) is with dominant diagonal, i.e.

$$|b| \geq |a| + |c|$$

also we must suppose $ac \neq 0$. (In the other case the system can be solved immediately).

In this case [2] the sequence $\{\alpha_i\}$ has the following properties:

a) $|\alpha_i| \leq 1$  for $i = 0, 1, 2, ...$

b) $\alpha_i \to \alpha$  for $i \to \infty$

where $a\alpha^2 + b\alpha + c = 0$.

Let us now study the speed of convergence of the sequence $\{\alpha_i\}$. For this aim we consider the following iteration process

(2)                    $$\alpha_i = -c(\alpha_{i-1}a + b)^{-1} \quad (i = 1, 2, 3...)$$

with $\alpha_i = \alpha + \epsilon_i$ $(i = 0, 1, 2...)$. Then from (2) we obtain that $\epsilon_i$ satisfy the difference equation

(3)                    $$\epsilon_i = \frac{\alpha q \epsilon_{i-1}}{\alpha - q\epsilon_{i-1}}, \quad q \equiv a\alpha^2/c.$$

To solve (3), first we are to study the case, when

(4)                    $$|a\alpha^2/c| = 1$$

i.e. $a\alpha^2/c = \pm 1$. But from the last equality we obtain $\mp 1 + b\alpha/c + 1 = 0$ which is possible only when

(5)                    $$a\alpha^2/c = 1.$$

Here we obtain that $b\alpha + 2c = 0$, or $\alpha = -2c/b$. Putting this value for $\alpha$ in (5), we find that $4ac = b^2$, from here we obtained that $a$ and $c$ have one and the same sign. From the other hand

$$|a| + |c| \leq |b| = 2\sqrt{|a||c|} \leq |a| + |c|.$$

From this relation we receive $a = c$. Therefore $\alpha = \pm 1$. Now if (4) holds, i.e. $|q| = 1$, $a = c$, $\alpha = \pm 1$. By the same way in this case from (3) we receive the following two difference equations

$$\epsilon_i = \frac{\epsilon_{i-1}}{1 - \epsilon_{i-1}}; \quad \epsilon_i = \frac{\epsilon_{i-1}}{1 + \epsilon_{i-1}}$$

corresponding to the both values $\alpha = \mp 1$ with the following solutions

(6) $$\epsilon_i = -\frac{k}{1 + ki}; \quad \epsilon_i = \frac{k}{1 + ki}$$

correspondingly where $k$ is an arbitrary constant. When we use $\epsilon_0 = -\alpha$, from (6) we receive $k = -1$ or $k = 1$ respectively. Then (6) takes the form

(7) $$\epsilon_i = -\frac{1}{1 + i}; \quad \epsilon_i = \frac{1}{1 + i} \quad (i = 0, 1, 2, ...)$$

If $|a\alpha^2/c| \neq 1$, it is easy to see that $|q| = |a\alpha^2/c| < 1$. In this case we receive

(8) $$\epsilon_i = -\frac{q^i \alpha}{1 + q + q^2 + ... + q^i} \quad (i = 0, 1, 2...).$$

for the solution of (3).

The relation (7) shows that the bounds, obtained by J a n j g a v a [1] for $|\epsilon_i|$ when $q = 4|a| \cdot |c| b^{-2} = 1$ are exact, the relation (8) gives better values for $|\epsilon_i|$.

## 2. Block case

Let it be given the following nonsingular quasitridiagonal system

(9)
$$
\begin{aligned}
BX_1 + CX_2 & & = D_1 \\
AX_1 + BX_2 + CX_3 & & = D_2 \\
\cdots\cdots\cdots\cdots\cdots & & \cdots\cdots\cdots\cdots\cdots \\
\cdots\cdots\cdots\cdots\cdots & & \cdots\cdots\cdots\cdots\cdots \\
& AX_{p-1} + BX_p & = D_p
\end{aligned}
$$

where $A$, $B$, $C$ are $m \times m$ matrices, $X_i$ and $D_i$ are $m \times 1$ matrices.

By the first step on Sweep's method the first $p - 1$ equations of (9) are transformed in the following form

$$X_1 = L_1 X_2 + Y_1$$

$$X_2 = L_2 X_3 + Y_2$$

$$\dots\dots\dots\dots\dots$$

$$\dots\dots\dots\dots\dots$$

$$X_{p-1} = L_{p-1} X_p + Y_{p-1}$$

where

(10)        $$L_k = -(AL_{k-1} + B)^{-1}C, \quad Y_k = -(AL_{k-1} + B)^{-1}(D_k - AY_{k-1})$$

and $k = 1, 2, ..., p; \; L_0 = 0, \; Y_0 = 0.$

In [1] P.V. J a n j g a v a proved the following

**Theorem 1.** *Let $B$ be an nonsingular matrix and $\| \cdot \|$ be a matrix norm, with $\|I\| = 1$, where $I$ is identity matrix. If*

(11)                          $$q \equiv 4\left\| B^{-1}A \right\| \cdot \left\| B^{-1}C \right\| \le 1.$$

*Then the matrix sequence*

(12)                          $$L_k = (AL_{k-1} + N)^{-1}C$$

*for $k = 1, 2, 3, ...$ and $L_0 = 0$ is converging to the solution $L$ of the quadratic matrix equation*

$$AX^2 + BX + C = 0 \text{ with the estimate}$$

$$\|L_k - L\| \le \frac{q^k}{k+1}\|L\|.$$

If in (11) we have only inequality, i.e. $q < 1$, then

$$\|L_k - L\| \le \frac{q^k}{(1 + \sqrt{1 - q})^{2k}}\|L\|.$$

Now we continue the Janjgava's investigations for the block case. For this aim we rewrite (12) in the following form

$$L_k = (PL_{k-1} + I)^{-1}Q$$

with $P = B^{-1}A$, $Q = B^{-1}C$.

In this case (11) takes the following form

(13) $$q = 4\|P\| \cdot \|Q\| \leq 1.$$

From Theorem 1 it follows that if (13) is satisfied, then the map

$$FX = -(PX + I)^{-1}Q$$

is defined in $\|X\| \leq 2\|Q\|$ it has a fixed point $L$ in the domain $\|X\| \leq 2\|Q\|$. Now we are going to prove the following

**Theorem 2.** *If $q \in (0,1)$ and*

$$G = \{X = n \times n \, matrix : \|X\| \leq 2\|Q\|\}$$

*then the map $FX$ has no other fixed points in $G$ except $L$.*

Proof. First, we are going to show that $X \in G \Longrightarrow FX \in G$ (see [1]). Indeed, let $X \in G$, i.e. $\|X\| \leq 2\|Q\|$. Then

$$\|FX\| = \|(PX + I)^{-1}Q\| \leq \|(PZ + I)^{-1}\| \|Q\| \leq \frac{\|Q\|}{1 - \|P\| - \|X\|}$$

$$\leq \frac{\|Q\|}{1 - 2\|P\| \cdot \|X\|} \leq \frac{\|Q\|}{1 - 1/2} = 2\|Q\|$$

then $FX \in G$.

Now we finish the proof showing that $FX$ is a contacting map in $G$. Indeed, if $X, Y \in G$, then

$$\|FX - FY\| = \|(PX + I)^{-1}P(Y - X).(PY + I)^{-1}Q\|$$

$$\leq \|(PX + I)^{-1}\| \cdot \|P\| \cdot \|Y - X\| \cdot \|(PY + I)^{-1}\| \cdot \|Q\|$$

$$\leq \frac{\|P\|\|Q\|}{(1 - 2\|P\| \cdot \|Q\|)^2} \cdot \|X - Y\| = \frac{q/4}{(1 - q/2)^2} \cdot \|X - Y\| = \frac{q}{(2 - q)^2}\|X - Y\|.$$

Since $q/(2 - q)^2 \in (0,1)$, then we obtain that $FX$ is a contracting map in $G$, which is enough to see that $FX$ has a unique fixed point $L$ in $G$. ∎

Now with the suggestion $\det(PQ) \neq 0$, we will obtain the following analytical equation of the error

$$\epsilon_k = L_k - L.$$

For the iteration process

$$L_k = -(PL_{k-1} + I)^{-1}Q$$

we have

$$(PL_{k-1} + I)L_k + Q = 0$$

$$(PL + P_{\epsilon_{k-1}} + I)(L + \epsilon_k) + Q = 0$$

$$(PL + I + P_{\epsilon_{k-1}})\epsilon_k + P\epsilon_{k-1}L = 0$$

By an inductive argument it is easy to see that for each $k = 0, 1, 2, ...,$ we have $|\epsilon_k| \neq 0$. In this case we can put $\delta_k = \epsilon_k^{-1}$.

In the same way, we continue to investigate following differentce equation

$$L\delta_k + \delta_{k-1}(L + P^{-1}) + I = 0$$

for $\delta_k$, or

(14)                          $$L\delta_k - \delta_{k-1}M + I = 0$$

where

$$M = -(L + P^{-1}).$$

In order to solve (14) it is necessary to find the general solution of the homogeneous equation $L\delta_k - \delta_{k-1}M = 0$ and one special solution for the non-homogeneous.

The general solution of homogeneous equation is

$$\Delta_k = L^{-k}CM^k, \quad k = 0, 1, 2, ...$$

where $C$ is an arbitrary $n \times n$ matrix.

Now we are to find a particular solution of (14).

We search it in the following form

$$\eta_k = L^{-k-1}C_kM^k.$$

After in (14) putting $\delta_k = \xi_k$ we obtain

$$L^{-k}C_kM^k - L^{-k}C_{k-1}M^k + I = 0$$

$$C_k - C_{k-1} + L^kM^{-k} = 0.$$

From $C_0 = 0$, we find

$$C_k = -\sum_{s=1}^{k} L^s M^{-s}.$$

After this one particular solution for (14) is

$$\eta_k = -\sum_{s=1}^{k} L^{s-k-1} M^{k-s}.$$

For the general solution $\delta_k = \Delta_k + \eta_k$ we find

$$\delta_k = L^{-k} C M^k - \sum_{s=1}^{k} L^{s-k-1} M^{k-s}$$

i.e.

$$\epsilon_k = \left( L^{-k} C M^k - \sum_{s=1}^{k} L^{s-k-1} M^{k-s} \right)^{-1}$$

where $C$ is an arbitrary matrix. After determining the solution from the condition $\epsilon_0 = -L$, we find finally

$$\epsilon_k = -M^{-k} \left( I + \sum_{s=1}^{k} L^s M^{-s} \right) L^{k+1},$$

where $k = 0, 1, 2, \ldots$

The convergence property of the coefficients $L_k$ ensures a possibility for a modification of Sweep's method, described from J a n j g a v a [1], which is more economic with respect of capacities and of the time of calculating, compared with the usual Sweep's method.

This modification can be described as follows:

Let we are to solve the linear system (9) by the Sweep's method, for the coefficients $L_k$ converging to $L$. Moreover, let $\tau < p - 1$ is the smallest integer such that $L_k = L$ for $k > \tau$ with accuracy $\epsilon$. In this case we should not calculate

$$L_{\tau+1} = L_{\tau+2} = \ldots = L_{p-1} = L.$$

and for $k > \tau + 1$ we calculate by the formula $Y_k$

$$Y_k = (AL + B)^{-1}(D_k - AY_{k-1})$$

where again we have an economy of calculating time.

The above modification was experimented in the scalar and block case. The experimental results supported the theoretical scheme.

**References**

[1] P.V. J a n j g a v a. On a property of the coefficients of the Sweep's method. *Vic. Mat. Program*, Tbilisi, 1976, **5-13**, (Russian).

[2] B. S. S e n d o v, V.A. P o p o v. Numerical Methods, part 2, Sofia, 1978, (Bulgarian)

*Institute of Mathematics*
*Acad. G. Bonchev str. , bl.8*
*1113 Sofia*
*BULGARIA*