14/01/2025 16:36

(Article begins on next page)

# A nonlinearity lagging for the solution of nonlinear steady state reaction diffusion problems

Emanuele Galligani

*Dipartimento di Matematica Pura e Applicata "G. Vitali"*
*Università degli Studi di Modena e Reggio Emilia*
*Via Campi 213/b, I–41125, Modena, Italy*

## Abstract

This paper concerns with the analysis of the iterative procedure for the solution of a nonlinear reaction diffusion equation at the steady state in a two dimensional bounded domain supplemented by suitable boundary conditions. This procedure, called Lagged Diffusivity Functional Iteration (LDFI)–procedure, computes the solution by "lagging" the diffusion term. A model problem is considered and a finite difference discretization for that model problem is described. Furthermore, properties of the finite difference operator are proved. Then, sufficient conditions for the convergence of the LDFI–procedure are given. At each stage of the LDFI–procedure a weakly nonlinear algebraic system has to be solved and the simplified Newton–Arithmetic Mean (Newton–AM) method is used. This method is particularly well suited for implementation on parallel computers. Numerical studies show the efficiency, for different test functions, of the LDFI–procedure combined with the simplified Newton–AM method. Better results are obtained when in the reaction diffusion equation also a convection term is present.

*Keywords:* Nonlinear problems, lagging diffusivity, Newton method, Arithmetic Mean method.
*2000 MSC:* 65H10, 65N06, 65N22

## 1. Introduction

We consider a nonlinear steady state reaction diffusion equation where the diffusion coefficient and the rate of change due to a reaction depend on the solution. These two terms are denoted with $\sigma$ and $-g$, respectively.

When we use a finite difference discretization, this elliptic equation supplemented by a Dirichlet boundary condition, can be transcribed into a nonlinear system of algebraic equations.

We wish to compute a solution of this system of nonlinear equations with a common iterative procedure in which the nonlinear term, corresponding to the discretization of the diffusivity $\sigma$, may be evaluated at the previous iteration (see [30]).

This approach of *nonlinearity lagging* in the diffusivity term is denoted as *Lagged Diffusivity Functional Iteration* (LDFI) or Lagged Diffusivity Fixed Point iteration.

In Section 2, the LDFI–procedure is stated: the nonlinear difference system is solved via a sequence of systems of *weakly nonlinear* difference equations where only the term corresponding to $g$ is nonlinear.

Thus, the iterates of the LDFI–procedure are the approximate solutions of the weakly nonlinear systems computed with an *inner* iterative solver; a criterion for acceptability of these approximate solutions is given. A stopping rule for the LDFI–procedure is also given.

Since, a purpose here is to re-examine the LDFI–procedure for solving the system of nonlinear difference equations of elliptic type in the context of Parallel Computing, a simplified version of the Newton–Arithmetic Mean (Newton–AM) method ([9]) is the inner iterative solver used for the solution of the weakly nonlinear systems. This is a *two–stage* iterative method and is particularly suited for implementation on parallel computers ([7], [8]; see also [1], [2], [3], [32]). In Section 6 a description of the simplified Newton–AM method and a general result on the convergence are reported.

The purpose of Sections 4 and 5 is to analyse the convergence of the LDFI–procedure, considering the essential features of the system of nonlinear equations generated by a finite difference discretization of a reaction diffusion *model problem* described in Section 3.

It is important to define for this model problem the set, called $\mathcal{B}$, of all grid functions defined on the discretized domain which contains the solutions of the weakly nonlinear systems and all iterates of the LDFI–procedure.

For example, the model problem studied in [21] allows to present a helpful paradigm for proving the convergence of the LDFI–procedure.

In Section 4, we summarize some properties of finite difference operators defined in $\mathcal{B}$ that also imply the uniform monotonicity of the nonlinear mapping which defines the nonlinear difference system.

In Section 5, the convergence, to a solution of the original nonlinear system, of the sequence of iterates generated with the LDFI–procedure is proved under mild and reasonable assumptions imposed on the diffusivity $\sigma$ and on function $g$ using well known standard techniques.

In the section of the numerical experiments (Section 7) the behaviour of the inner–outer iterations of the procedure is examined. The effectiveness of the LDFI–procedure combined with the simplified Newton–AM method is highlighted, especially, for reaction diffusion problems where also the convection term is present.

## 2. The lagged diffusivity functional iteration procedure

Consider a *strongly nonlinear* system of algebraic equations of the form

$$\boldsymbol{F}(\boldsymbol{u}) \equiv A(\boldsymbol{u})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0}, \tag{1}$$

where $\boldsymbol{u} = (u_1, u_2, ..., u_n)^T$ is a vector in $\mathbb{R}^n$, $A(\boldsymbol{u})$ is a large $n \times n$ nonsingular matrix with a sparse structure, and $\boldsymbol{G}(\boldsymbol{u})$ is a continuously differentiable *diagonal mapping*, i.e., a nonlinear mapping whose $i$–th component $G_i$ is a function of only the $i$-th variable $u_i$ for $i = 1, ..., n$. $\boldsymbol{s}$ is a vector of $n$ components independent of $\boldsymbol{u}$. We assume that this system has a solution $\boldsymbol{u}^*$.

For solving system (1) the easiest and maybe the most common method is *to lag* part of the nonlinear terms in (1) generating an iterative procedure denoted as Lagged Diffusivity Functional Iteration.

With this iterative procedure the nonlinear system (1) can be solved via a sequence of systems of *weakly nonlinear* equations where the nonlinear term is $\boldsymbol{G}$.

Specifically, given a sequence of positive numbers $\{\varepsilon_\nu\}$ such that $\varepsilon_\nu \to 0$ as $\nu \to \infty$ and an initial estimate $\boldsymbol{u}^{(0)}$ of the solution $\boldsymbol{u}^*$ of the system (1), we generate a sequence of iterates $\{\boldsymbol{u}^{(\nu)}\}$, $\nu = 0, 1, 2, ...$, with the following rule for the transition from a current iteration $\boldsymbol{u}^{(\nu)}$ to the new iterate $\boldsymbol{u}^{(\nu+1)}$:

- Find an approximate solution $\boldsymbol{u}^{(\nu+1)}$ of the nonlinear system

$$\boldsymbol{F}_\nu(\boldsymbol{u}) \equiv A(\boldsymbol{u}^{(\nu)})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0}, \tag{2}$$

  with the criterion for acceptability of the solution

$$\|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \le \varepsilon_{\nu+1}. \tag{3}$$

Then, the LDFI–procedure is composed by a nonlinear *outer* iteration that generates the sequence $\{\boldsymbol{u}^{(\nu)}\}$ and by an *inner* iterative solver of the weakly nonlinear system (2). This solver must be particularly well suited for implementation on parallel computers.

The termination criterion for the outer iteration is provided by the following two inequalities

$$\|\boldsymbol{u}^{(\nu+1)} - \boldsymbol{u}^{(\nu)}\| \le \tau_1,$$

$$\tag{4}$$

$$\|\boldsymbol{F}(\boldsymbol{u}^{(\nu+1)})\| \le \tau_2,$$

where $\tau_1$ and $\tau_2$ are prespecified tolerances ([10, Theor. 3]).

For the inner iterative solver, we use the following rule

$$\varepsilon_{\nu+1} = 0.5\varepsilon_\nu, \qquad \nu = 1, 2, ...$$

with $\varepsilon_1 = 0.1\,\|\boldsymbol{F}(\boldsymbol{u}^{(0)})\|$.

## 3. A model problem

Consider a nonlinear steady state reaction diffusion convection equation of the form

$$-\mathrm{div}(\sigma(\boldsymbol{x}, \varphi)\nabla\varphi) + \tilde{\boldsymbol{v}} \cdot \nabla\varphi + \alpha(\boldsymbol{x})\varphi + g(\boldsymbol{x}, \varphi) = s(\boldsymbol{x}) \qquad \boldsymbol{x} \in \Omega, \tag{5}$$

where $\varphi = \varphi(\boldsymbol{x})$ is the density function at the point $\boldsymbol{x}$ of a diffusion medium $\Omega$, $\sigma = \sigma(\boldsymbol{x}, \varphi) > 0$ is the diffusion coefficient or diffusivity and is dependent on the solution $\varphi$, $\alpha = \alpha(\boldsymbol{x}) \ge 0$ is the absorption term, $\tilde{\boldsymbol{v}}$ is the velocity vector, $-g(\boldsymbol{x}, \varphi)$ is the rate of change due to a reaction and $s(\boldsymbol{x})$ is the source term.

Here, the vector $\tilde{\boldsymbol{v}} = (\tilde{v}_1, \tilde{v}_2)^T$ is assumed to be constant.

In the equation (5) the convection term $\tilde{\boldsymbol{v}} \cdot \nabla \varphi$ has been taken into account; however, we will consider only convection–not dominated problems.

We consider that equation (5) is subject by homogeneous Dirichlet boundary condition on the contour $\Gamma$ of $\Omega$:

$$\varphi(\boldsymbol{x}) = 0 \qquad \boldsymbol{x} \in \Gamma. \tag{6}$$

The functions $\sigma(\boldsymbol{x}, \varphi)$, $\alpha(\boldsymbol{x})$, $s(\boldsymbol{x})$ and $g(\boldsymbol{x}, \varphi)$ are assumed to satisfy the following "smoothness" conditions:

**(i)** the functions $\sigma(\boldsymbol{x}, \varphi)$ and $g(\boldsymbol{x}, \varphi)$ are continuously differentiable in $\boldsymbol{x}$ and continuous in $\varphi$; the functions $\alpha(\boldsymbol{x})$ and $s(\boldsymbol{x})$ (the "source term") are continuous in $\boldsymbol{x}$;

**(ii)** there exist two positive constants $\sigma_{\min}$ and $\sigma_{\max}$ such that

$$0 < \sigma_{\min} \leq \sigma(\boldsymbol{x}, \varphi) \leq \sigma_{\max},$$

uniformly in $\boldsymbol{x}$ and $\varphi$; in addition, $\alpha(\boldsymbol{x}) \geq 0$;

**(iii)** for fixed $\boldsymbol{x} \in \Omega$, the function $\sigma(\boldsymbol{x}, \varphi)$ satisfies Lipschitz condition in $\varphi$ with constant $\Lambda$ (uniformly in $\boldsymbol{x}$), $\Lambda > 0$;

**(iv)** for a fixed $\boldsymbol{x} \in \Omega$, the function $g(\boldsymbol{x}, \varphi)$ is a uniformly monotone mapping ([26, p. 141]) in $\varphi$ with constant $c > 0$ (uniformly in $\boldsymbol{x}$) and is continuously differentiable in $\varphi$.

We assume that the problem (5)–(6) has an isolated solution.

There exist various techniques for discretizing the problem (5)–(6). Using the Taylor series approach, equation (5) will be solved with the following standard finite difference scheme.

We consider $\Omega$ a rectangular domain ($\boldsymbol{x} \equiv (x, y)^T$) with boundary $\Gamma$ and we superimpose on $\Omega \cup \Gamma$ a grid of points $\Omega_h \cup \Gamma_h$; the set of the internal points $\Omega_h$ of the grid are the mesh points $(x_i, y_j)$, for $i = 1, ..., N$ and $j = 1, ..., M$, with uniform mesh size $h$ along $x$ and $y$ directions respectively, i.e. $x_{i+1} = x_i + h$ and $y_{j+1} = y_j + h$ for $i = 0, ..., N$, $j = 0, ..., M$.

Furthermore, at the mesh points of $\Omega \cup \Gamma$, $(x_i, y_j)$, for $i = 0, ..., N + 1$ and $j = 0, ..., M + 1$, the solution $\varphi(x_i, y_j)$ is approximated by a *grid function $u_{ij}$* defined on $\Omega_h \cup \Gamma_h$ and vanishing on $\Gamma_h$.

In order to approximate partial derivatives in (5) we shall make use of difference quotients of grid functions. The forward, backward and centered

difference quotients with respect to $x$ and to $y$ of the grid function $u_{ij}$ at the mesh point $(x_i, y_j)$, are, respectively:

$$\Delta_x u_{ij} = \frac{u_{i+1j} - u_{ij}}{h}, \qquad \Delta_y u_{ij} = \frac{u_{ij+1} - u_{ij}}{h},$$

$$\nabla_x u_{ij} = \frac{u_{ij} - u_{i-1j}}{h}, \qquad \nabla_y u_{ij} = \frac{u_{ij} - u_{ij-1}}{h},$$

$$\delta_x u_{ij} = \frac{1}{2}(\Delta_x u_{ij} + \nabla_x u_{ij}), \qquad \delta_y u_{ij} = \frac{1}{2}(\Delta_y u_{ij} + \nabla_y u_{ij}),$$

while the centered second difference quotient with respect to $x$ and to $y$ can be written

$$\delta_x^2 u_{ij} = \nabla_x \Delta_x u_{ij} = \Delta_x \nabla_x u_{ij}, \qquad\qquad \delta_y^2 u_{ij} = \nabla_y \Delta_y u_{ij} = \Delta_y \nabla_y u_{ij}.$$

This notation was introduced in [6].

Providing a discretization error $O(h^2)$, the finite difference approximation of (5) in $(x_i, y_j)$ is given by

$$-\Delta_x \left( \sigma(x_i, y_j, u_{ij}) \nabla_x u_{ij} \right) - \Delta_y \left( \sigma(x_i, y_j, u_{ij}) \nabla_y u_{ij} \right) + \tilde{v}_1 \delta_x u_{ij} + \tilde{v}_2 \delta_y u_{ij} +$$
$$+ \alpha(x_i, y_j) u_{ij} + g(x_i, y_j, u_{ij}) = s(x_i, y_j),$$

that yields to

$$-(B_{ij} + \tilde{B}_{ij}) u_{ij-1} - (L_{ij} + \tilde{L}_{ij}) u_{i-1j} + (D_{ij} + \tilde{D}_{ij}) u_{ij} - \qquad (7)$$
$$-(R_{ij} + \tilde{R}_{ij}) u_{i+1j} - (T_{ij} + \tilde{T}_{ij}) u_{ij+1} + g(x_i, y_j, u_{ij}) - s(x_i, y_j) = 0,$$

where the coefficients are: for $i = 1, ..., N$ and $j = 1, ..., M$

$$L_{ij} \equiv L_{ij}(\boldsymbol{u}) = \frac{1}{h^2} \sigma(x_i, y_j, u_{ij}), \qquad B_{ij} \equiv B_{ij}(\boldsymbol{u}) = \frac{1}{h^2} \sigma(x_i, y_j, u_{ij}),$$

$$R_{ij} \equiv R_{ij}(\boldsymbol{u}) = \frac{1}{h^2} \sigma(x_{i+1}, y_j, u_{i+1j}), \qquad T_{ij} \equiv T_{ij}(\boldsymbol{u}) = \frac{1}{h^2} \sigma(x_i, y_{j+1}, u_{ij+1}),$$

$$\tilde{L}_{ij} = \frac{\tilde{v}_1}{2h}, \qquad \tilde{B}_{ij} = \frac{\tilde{v}_2}{2h}, \qquad\qquad (8)$$

$$\tilde{R}_{ij} = -\frac{\tilde{v}_1}{2h}, \qquad \tilde{T}_{ij} = -\frac{\tilde{v}_2}{2h},$$

$$D_{ij} \equiv D_{ij}(\boldsymbol{u}) = B_{ij} + L_{ij} + R_{ij} + T_{ij}, \qquad \tilde{D}_{ij} = \alpha(x_i, y_j).$$

We denote the mesh points $P_k = (x_i, y_j)$, $(i = 1, ..., N, j = 1, ..., M)$ and we order the points $P_k$ in lexicographic order: $k = (j - 1) \times N + i$. We set

$n = N \times M$, and we denote the vector solution $\boldsymbol{u}$ whose components are the values of the grid function at the internal mesh points

$$\boldsymbol{u} = (u_1, ..., u_n)^T \equiv (u_{11}, ..., u_{N1}, u_{12}, ..., u_{N2}, ..., u_{1M}, ..., u_{NM})^T.$$

Then, formula (7) for $i = 1, ..., N$, $j = 1, ..., M$, can be written as formula (1)

$$\boldsymbol{F}(\boldsymbol{u}) \equiv A(\boldsymbol{u})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0}, \tag{9}$$

where the matrix $A(\boldsymbol{u})$ of order $n$ has a block tridiagonal form.

The $M$ diagonal blocks are tridiagonal matrices of order $N$ with diagonal elements $a_{kk}(\boldsymbol{u}) = D_{ij} + \tilde{D}_{ij}$, the sub– and super–diagonal elements are $a_{k-1k}(\boldsymbol{u}) = -(L_{ij} + \tilde{L}_{ij})$ and $a_{kk+1}(\boldsymbol{u}) = -(R_{ij} + \tilde{R}_{ij})$ respectively, $i = 1, ..., N$, $j = 1, ..., M$ and $k = (j-1) \times N + i$ (here $L_{1j}, \tilde{L}_{1j}, R_{Nj}$ and $\tilde{R}_{Nj}$, $j = 1, ..., M$, are the coefficients of the solution computed in mesh points $\Gamma_h$).
The sub– and super–diagonal blocks are diagonal matrices of order $N$ with elements $a_{k-Nk}(\boldsymbol{u}) = -(B_{ij} + \tilde{B}_{ij})$ and $a_{kk+N}(\boldsymbol{u}) = -(T_{ij} + \tilde{T}_{ij})$ respectively, $i = 1, ..., N$, $j = 1, ..., M$ and $k = (j-1) \times N + i$ (here $B_{i1}, \tilde{B}_{i1}, T_{iM}$ and $\tilde{T}_{iM}$, $i = 1, ..., N$, are the coefficients of the solution computed in mesh points $\Gamma_h$).
Providing that the mesh spacing $h$ is sufficiently small, i.e.

$$h < \min\left\{\frac{2\sigma_{\min}}{|\tilde{v}_1|}, \frac{2\sigma_{\min}}{|\tilde{v}_2|}\right\},$$

the matrix $A(\boldsymbol{u})$ is strictly ($\alpha(x, y) > 0$) or irreducibly ($\alpha(x, y) = 0$) diagonally dominant ([31, p. 23]) and has positive diagonal elements, $a_{kk}(\boldsymbol{u}) > 0$ and nonpositive off diagonal elements $a_{kl}(\boldsymbol{u}) \leq 0, k \neq l$, with $k, l = 1, ..., n$; therefore $A(\boldsymbol{u})$ is an M–matrix ([31, p. 91]).
In the case of reaction diffusion equation ($\tilde{\boldsymbol{v}} = \boldsymbol{0}$), the matrix $A(\boldsymbol{u})$ is also symmetric; then $A(\boldsymbol{u})$ is symmetric positive definite (Stieltjes matrix [31, p. 91]).
In the following, we may consider the matrix $A(\boldsymbol{u})$ as

$$A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{A} + \tilde{D},$$

where $A_1(\boldsymbol{u})$ and $\tilde{A}$ are the block tridiagonal matrices containing the elements $\{B_{ij}, L_{ij}, D_{ij}, R_{ij}, T_{ij}\}$ and $\{\tilde{B}_{ij}, \tilde{L}_{ij}, \tilde{R}_{ij}, \tilde{T}_{ij}\}$ respectively, while the matrix $\tilde{D}$ is a diagonal matrix whose diagonal entries are $\{\tilde{D}_{ij}\}$.
Furthermore, $\boldsymbol{s} \in \mathbb{R}^n$ is a vector whose components are the values of the source term $s(x, y)$ at the mesh points; the nonlinear mapping $\boldsymbol{G}(\boldsymbol{u})$ has

components $G_k(\boldsymbol{u}) = g(x_i, y_j, u_k)$, $i = 1, ..., N$, $j = 1, ..., M$ and $k = (j - 1) \times N + i$. We observe, that $G_k(\boldsymbol{u})$, the $k$–th component of $\boldsymbol{G}(\boldsymbol{u})$, with respect to the variable $\boldsymbol{u}$, depends of only the $k$–th component $u_k$, for $k = 1, ..., n$; in this case $\boldsymbol{G}$ is a diagonal mapping.

We observe that the right hand side of (9) is the null vector since we have the homogeneous Dirichlet condition (6) in $\Gamma$.

For grid functions $\{u_{ij}\}$ and $\{v_{ij}\}$ of this type the discrete $l_2(\Omega_h)$ inner product and norm are defined by the formulas

$$< \boldsymbol{u}, \boldsymbol{v} > = h^2 \sum_{i=1}^{N} \sum_{j=1}^{M} u_{ij} v_{ij},$$

(10)

$$\|\boldsymbol{u}\|_h = (h^2 \sum_{i=1}^{N} \sum_{j=1}^{M} |u_{ij}|^2)^{1/2} = (< \boldsymbol{u}, \boldsymbol{u} >)^{1/2},$$

respectively.

We say that the grid functions $\{u_{ij}\}$ defined on $\Omega_h \cup \Gamma_h$ and vanishing on $\Gamma_h$ satisfy **Property A** if they are uniformly bounded and have uniformly bounded backward difference quotients $\nabla_x u_{ij}$ and $\nabla_y u_{ij}$ at each mesh point $(x_i, y_j)$ of $\Omega_h \cup \Gamma_h$. The set of all grid functions $\{u_{ij}\}$ which satisfy Property A is denoted by $\mathcal{B}$. Thus, $\mathcal{B}$ is the set of grid functions $\{u_{ij}\}$ for which there exist some positive constants $\rho$ and $\beta$ such that

$$\|\boldsymbol{u}\|_h \leq \rho,$$

(11)

$$|\nabla_x u_{ij}| \leq \beta \quad \text{and} \quad |\nabla_y u_{ij}| \leq \beta.$$

(12)

The constant $\rho$ is independent of $h$; also the constant $\beta$ is independent of $h$ but it depends on $\|\boldsymbol{G}(\boldsymbol{u}) + \boldsymbol{s}\|_h$.

We assume that the system (9) has at least one solution $\boldsymbol{u}^*$ in $\mathcal{B}$ with $|\nabla_x u_{ij}^*| \leq \beta$ and $|\nabla_y u_{ij}^*| \leq \beta$.

(See, i.e., [21], where a proof of the existence of such a solution $\boldsymbol{u}^*$ of (9) in $\mathcal{B}$ is given; also a condition for which $\boldsymbol{u}^*$ is unique in $\mathcal{B}$ has been obtained)

## 4. Some properties of finite difference operators

In the following we summarize some properties of finite difference operators when $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ and $\alpha(x, y) = 0$ in (5).

Here we denote

$$\sigma_{lp}(u_{ij}) \equiv \sigma(x_l, y_p, u_{ij}), \qquad g_{lp}(u_{ij}) \equiv g(x_l, y_p, u_{ij}).$$

**Lemma 1.** Let $\{u_{ij}\}$, $\{v_{ij}\}$, $\{w_{ij}\}$ be three grid functions defined at the mesh points $(x_i, y_j)$ of a grid $\Omega_h \cup \Gamma_h$, $i = 0, ..., N+1$, $j = 0, ..., M+1$ which are zero on $\Gamma_h$. Suppose the coefficients in (8), $L_{ij}$, $R_{ij}$, $B_{ij}$ and $T_{ij}$, are dependent on the grid function $w_{ij}$, then,

$$\sum_{i=1}^{N} [L_{ij}(u_{ij} - u_{i-1j}) - R_{ij}(u_{i+1j} - u_{ij})] v_{ij} = \tag{13}$$

$$= \sum_{i=1}^{N} \left[ \frac{1}{h^2} \sigma_{ij}(w_{ij})(u_{ij} - u_{i-1j})(v_{ij} - v_{i-1j}) \right] + \frac{1}{h^2} \sigma_{N+1j}(w_{N+1j}) u_{Nj} v_{Nj},$$

and

$$\sum_{j=1}^{M} [B_{ij}(u_{ij} - u_{ij-1}) - T_{ij}(u_{ij+1} - u_{ij})] v_{ij} = \tag{14}$$

$$= \sum_{j=1}^{M} \left[ \frac{1}{h^2} \sigma_{ij}(w_{ij})(u_{ij} - u_{ij-1})(v_{ij} - v_{ij-1}) \right] + \frac{1}{h^2} \sigma_{iM+1}(w_{iM+1}) u_{iM} v_{iM}.$$

**Proof.** We prove formula (13). We have[1]

$$\sum_{i=1}^{N} [-L_{ij} u_{i-1j} + (L_{ij} + R_{ij}) u_{ij} - R_{ij} u_{i+1j}] v_{ij} =$$

$$= \sum_{i=1}^{N} [L_{ij}(u_{ij} - u_{i-1j}) - R_{ij}(u_{i+1j} - u_{ij})] v_{ij}$$

$$= \sum_{i=1}^{N} \left[ \frac{1}{h^2} \sigma(w_{ij})(u_{ij} - u_{i-1j}) - \frac{1}{h^2} \sigma(w_{i+1j})(u_{i+1j} - u_{ij}) \right] v_{ij}$$

$$= \frac{1}{h^2} \sigma(w_{1j})(u_{1j} - u_{0j}) v_{1j} - \frac{1}{h^2} \sigma(w_{2j})(u_{2j} - u_{1j}) v_{1j} +$$

$$+ \frac{1}{h^2} \sigma(w_{2j})(u_{2j} - u_{1j}) v_{2j} - \frac{1}{h^2} \sigma(w_{3j})(u_{3j} - u_{2j}) v_{2j} +$$

---

[1]For simplicity of notation, in the proofs of the Lemmas 1, 2 and 3 we omit the indexes of the coordinates $x$ and $y$ in the expression of the function $\sigma$.

$$+\frac{1}{h^2}\sigma(w_{3j})(u_{3j}-u_{2j})v_{3j}-\frac{1}{h^2}\sigma(w_{4j})(u_{4j}-u_{3j})v_{3j}+...$$

$$...+\frac{1}{h^2}\sigma(w_{N-1j})(u_{N-1j}-u_{N-2j})v_{N-1j}-$$

$$-\frac{1}{h^2}\sigma(w_{Nj})(u_{Nj}-u_{N-1j})v_{N-1j}+\frac{1}{h^2}\sigma(w_{Nj})(u_{Nj}-u_{N-1j})v_{Nj}-$$

$$-\frac{1}{h^2}\sigma(w_{N+1j})(u_{N+1j}-u_{Nj})v_{Nj},$$

then, since $v_{0j}=0$ for (6), the expression of the right hand side becomes

$$\frac{1}{h^2}\sigma(w_{1j})(u_{1j}-u_{0j})(v_{1j}-v_{0j})+\frac{1}{h^2}\sigma(w_{2j})(u_{2j}-u_{1j})(v_{2j}-v_{1j})+$$

$$+\frac{1}{h^2}\sigma(w_{3j})(u_{3j}-u_{2j})(v_{3j}-v_{2j})+...+\frac{1}{h^2}\sigma(w_{Nj})(u_{Nj}-u_{N-1j})\times$$

$$\times(v_{Nj}-v_{N-1j})-\frac{1}{h^2}\sigma(w_{N+1j})(u_{N+1j}-u_{Nj})v_{Nj},$$

and by $u_{N+1j}=0$, we have formula (13). Similarly, we obtain formula (14). $\sharp$

We remark that from the right hand side of (13) and (14) we can swap $u_{ij}$ with $v_{ij}$ and we obtain

$$\sum_{i=1}^{N}\left[L_{ij}(u_{ij}-u_{i-1j})-R_{ij}(u_{i+1j}-u_{ij})\right]v_{ij}=$$

$$=\sum_{i=1}^{N}\left[L_{ij}(v_{ij}-v_{i-1j})-R_{ij}(v_{i+1j}-v_{ij})\right]u_{ij},$$

and

$$\sum_{j=1}^{M}\left[B_{ij}(u_{ij}-u_{ij-1})-T_{ij}(u_{ij+1}-u_{ij})\right]v_{ij}=$$

$$=\sum_{j=1}^{M}\left[B_{ij}(v_{ij}-v_{ij-1})-T_{ij}(v_{ij-1}-v_{ij})\right]u_{ij}.$$

**Lemma 2.** Let $\{u_{ij}\}$, $\{v_{ij}\}$, $\{w_{ij}\}$ be three grid functions defined at the mesh points $(x_i,y_j)$ of a grid $\Omega_h\cup\Gamma_h$, $i=0,...,N+1$, $j=0,...,M+1$ which are zero on $\Gamma_h$.

Then, we have the following expression for the discrete $l_2(\Omega_h)$ inner product of the vectors $A(\boldsymbol{w})\boldsymbol{u}$ and $\boldsymbol{v}$ where the $n \times n$ matrix $A(\boldsymbol{w})$ is the one in (9), replacing $\boldsymbol{u}$ with $\boldsymbol{w}$:

$$
\begin{aligned}
< A(\boldsymbol{w})\boldsymbol{u}, \boldsymbol{v} > \ &= \ h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \sigma_{ij}(w_{ij})(\nabla_x u_{ij} \nabla_x v_{ij} + \nabla_y u_{ij} \nabla_y v_{ij}) + \quad (15) \\
&\quad + \sum_{j=1}^{M} \sigma_{N+1j}(w_{N+1j})u_{Nj}v_{Nj} + \sum_{i=1}^{N} \sigma_{iM+1}(w_{iM+1})u_{iM}v_{iM}.
\end{aligned}
$$

**Proof.** Suppose the coefficients in (8), $L_{ij}$, $R_{ij}$, $B_{ij}$ and $T_{ij}$, are functions of the grid function $w_{ij}$, we have,

$$
\begin{aligned}
< A(\boldsymbol{w})\boldsymbol{u}, \boldsymbol{v} > \ &= \ h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} [ -B_{ij} u_{ij-1} - L_{ij} u_{i-1j} + \\
&\qquad + (B_{ij} + L_{ij} + R_{ij} + T_{ij})u_{ij} - R_{ij} u_{i+1j} - T_{ij} u_{ij+1}] \, v_{ij} \\
&= \ h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} [ -B_{ij}(u_{ij} - u_{ij-1}) - L_{ij}(u_{ij} - u_{i-1j}) - \\
&\qquad - R_{ij}(u_{i+1j} - u_{ij}) - T_{ij}(u_{ij+1} - u_{ij})] \, v_{ij}.
\end{aligned}
$$

Using formulae (13) and (14) and keeping into account of the inner product in $l_2(\Omega_h)$, we have

$$
\begin{aligned}
< A(\boldsymbol{w})\boldsymbol{u}, \boldsymbol{v} > \ &= \ h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \frac{1}{h^2} \Big[ \sigma(w_{ij})(u_{ij} - u_{i-1j})(v_{ij} - v_{i-1j}) + \\
&\qquad + \sigma(w_{ij})(u_{ij} - u_{ij-1})(v_{ij} - v_{ij-1}) \Big] \ + \\
&\quad + \sum_{j=1}^{M} \sigma(w_{N+1j})u_{Nj}v_{Nj} + \sum_{i=1}^{N} \sigma(w_{iM+1})u_{iM}v_{iM} \\
&= \ h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \frac{1}{h^2} \Big[ \sigma(w_{ij}) \frac{u_{ij} - u_{i-1j}}{h} \frac{v_{ij} - v_{i-1j}}{h} h^2 + \\
&\qquad + \sigma(w_{ij}) \frac{u_{ij} - u_{ij-1}}{h} \frac{v_{ij} - v_{ij-1}}{h} h^2 \Big] \ + \\
&\quad + \sum_{j=1}^{M} \sigma(w_{N+1j})u_{Nj}v_{Nj} + \sum_{i=1}^{N} \sigma(w_{iM+1})u_{iM}v_{iM}.
\end{aligned}
$$

11

Then we have formula (15).     ♯

We remark that while the grid function $\{u_{ij}\}$ is defined on the entire mesh region $\Omega_h \cup \Gamma_h$, the vector $\boldsymbol{u} \in \mathbb{R}^n$ represents the grid function $\{u_{ij}\}$ defined only on the interior mesh points $\Omega_h$, $i = 1, ..., N$, $j = 1, ..., M$.

Moreover, we observe that formula (15) implies

$$< A(\boldsymbol{w})\boldsymbol{u}, \boldsymbol{v} > \;=\; < \boldsymbol{u}, A(\boldsymbol{w})\boldsymbol{v} > .$$

**Lemma 3.** Let $\{u_{ij}\}$, $\{v_{ij}\}$, $\{w_{ij}\}$ be three grid functions defined at the mesh points $(x_i, y_j)$ of a grid $\Omega_h \cup \Gamma_h$, $i = 0, ..., N+1$, $j = 0, ..., M+1$ which are zero on $\Gamma_h$.

Let $A(\boldsymbol{u})$ the matrix $n \times n$ in (9) and let $A(\boldsymbol{w})$ the matrix $n \times n$ in (9) with $\boldsymbol{u}$ replaced by the vector $\boldsymbol{w}$.

Then, if $\boldsymbol{u}$, $\boldsymbol{v}$ and $\boldsymbol{w}$ belong to $\mathcal{B}$, we have the following inequality

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \;\leq\; \frac{h^2 \Lambda \beta \phi}{2} \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + \right.$$

$$\left. + |\nabla_y(u_{ij} - v_{ij})|^2 \right] + \frac{\Lambda \beta}{\phi} \|\boldsymbol{u} - \boldsymbol{w}\|_h^2, \tag{16}$$

where $\Lambda > 0$ is the Lipschitz constant of condition (iii), $\beta > 0$ is a constant for which $|\nabla_x v_{ij}| \leq \beta$ and $|\nabla_y v_{ij}| \leq \beta$, and $\phi$ is an arbitrary positive number.

**Proof.** By using formula (15) in Lemma 2, we can write

$$< (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >=< A(\boldsymbol{u})\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > - < A(\boldsymbol{w})\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >=$$

$$= \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ h^2 \sigma(u_{ij})(\nabla_x v_{ij} \nabla_x(u_{ij} - v_{ij}) + \nabla_y v_{ij} \nabla_y(u_{ij} - v_{ij})) \right] +$$

$$+ \sum_{j=1}^{M} \sigma(u_{N+1j})v_{Nj}(u_{Nj} - v_{Nj}) + \sum_{i=1}^{N} \sigma(u_{iM+1})v_{iM}(u_{iM} - v_{iM}) -$$

$$- \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ h^2 \sigma(w_{ij})(\nabla_x v_{ij} \nabla_x(u_{ij} - v_{ij}) + \nabla_y v_{ij} \nabla_y(u_{ij} - v_{ij})) \right] -$$

$$- \sum_{j=1}^{M} \sigma(w_{N+1j})v_{Nj}(u_{Nj} - v_{Nj}) - \sum_{i=1}^{N} \sigma(w_{iM+1})v_{iM}(u_{iM} - v_{iM})$$

12

$$= \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ h^2 (\sigma(u_{ij}) - \sigma(w_{ij})) \left( \nabla_x v_{ij} \nabla_x (u_{ij} - v_{ij}) + \right. \right.$$

$$\left. + \nabla_y v_{ij} \nabla_y (u_{ij} - v_{ij}) \right)] + \sum_{j=1}^{M} (\sigma(u_{N+1j}) - \sigma(w_{N+1j})) v_{Nj} (u_{Nj} - v_{Nj}) +$$

$$+ \sum_{i=1}^{N} (\sigma(u_{iM+1}) - \sigma(w_{iM+1})) v_{iM} (u_{iM} - v_{iM}).$$

Now, we have that the term $| < (A(\boldsymbol{u}) - A(\boldsymbol{w})) \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > |$ is equal to the absolute value of the last expression. Then,

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w})) \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \le \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ h^2 |\sigma(u_{ij}) - \sigma(w_{ij})| \left( |\nabla_x v_{ij}| \times \right. \right.$$

$$\left. \times |\nabla_x (u_{ij} - v_{ij})| + |\nabla_y v_{ij}| |\nabla_y (u_{ij} - v_{ij})| \right)] +$$

$$+ \sum_{j=1}^{M} |\sigma(u_{N+1j}) - \sigma(w_{N+1j})| |v_{Nj}| |u_{Nj} - v_{Nj}| +$$

$$+ \sum_{i=1}^{N} |\sigma(u_{iM+1}) - \sigma(w_{iM+1})| |v_{iM}| |u_{iM} - v_{iM}|. \tag{17}$$

The assumption (iii) implies that, for given grid functions $\{u_{ij}\}$, $\{w_{ij}\}$ belonging to $\Omega_h \cup \Gamma_h$ there exists a positive constant $\Lambda$ such that for all $i = 1, ..., N+1$ and $j = 1, ..., M+1$

$$|\sigma(u_{ij}) - \sigma(w_{ij})| \le \Lambda |u_{ij} - w_{ij}|. \tag{18}$$

The constant $\Lambda$ is independent of $h$.

Furthermore, Property A assures that there exists a constant $\beta > 0$ such that inequality (12) holds

$$|\nabla_x v_{ij}| \le \beta \quad \text{and} \quad |\nabla_y v_{ij}| \le \beta,$$

for all $i = 1, ..., N+1$ and $j = 1, ..., M+1$ and all grid function $\{v_{ij}\}$ belonging to $\mathcal{B}$. The constant $\beta$ is independent of $h$.

Now, if we apply inequalities (18) and (12) into the expression in (17) we obtain

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w})) \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \le \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ h^2 \Lambda \beta |u_{ij} - w_{ij}| \left( |\nabla_x (u_{ij} - v_{ij})| + \right. \right.$$

13

$$+|\nabla_y(u_{ij} - v_{ij})|)] + \Lambda \sum_{j=1}^{M} |u_{N+1j} - w_{N+1j}| \, |v_{Nj}| \, |u_{Nj} - v_{Nj}| +$$

$$+\Lambda \sum_{i=1}^{N} |u_{iM+1} - w_{iM+1}| \, |v_{iM}| \, |u_{iM} - v_{iM}|.$$

Since the grid functions belonging to $\mathcal{B}$ are bounded and are equal to zero at the points of the boundary, we obtain

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \quad \leq \quad h^2 \Lambda \beta \sum_{j=1}^{M} \sum_{i=1}^{N} [|u_{ij} - w_{ij}| \, |\nabla_x(u_{ij} - v_{ij})| +$$

$$+ \, |u_{ij} - w_{ij}| \, |\nabla_y(u_{ij} - v_{ij})|] \, .$$

The last expression can be written

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \quad \leq \quad h^2 \Lambda \beta \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ \frac{|u_{ij} - w_{ij}|}{\sqrt{\phi}} \, |\nabla_x(u_{ij} - v_{ij})| \sqrt{\phi} + \right.$$

$$\left. + \, \frac{|u_{ij} - w_{ij}|}{\sqrt{\phi}} \, |\nabla_y(u_{ij} - v_{ij})| \sqrt{\phi} \right].$$

Using a well known technical trick, we have

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \leq h^2 \Lambda \beta \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ \frac{|u_{ij} - w_{ij}|^2}{2\phi} + \right.$$

$$\left. +|\nabla_x(u_{ij} - v_{ij})|^2 \frac{\phi}{2} + \frac{|u_{ij} - w_{ij}|^2}{2\phi} + |\nabla_y(u_{ij} - v_{ij})|^2 \frac{\phi}{2} \right], \quad (19)$$

and we obtain formula (16).  ♯

As consequence of Lemmas 1, 2 and 3, it is possible to prove that the mapping $\boldsymbol{F}(\boldsymbol{u})$ is uniformly monotone in $\mathcal{B}$ if the condition

$$\gamma \equiv \frac{\Lambda^2 \beta^2}{2\sigma_{\min} c} < 1,$$

is satisfied (see [21] (and [12])). Thus, from Hadamard Theorem ([17]), the nonlinear system (9) has a unique solution (e.g. [26, p. 143]).

Note that the two hypotheses that $\boldsymbol{F}(\boldsymbol{u})$ is Lipschitz–continuous and uniformly monotone on $\mathbb{R}^n$ are sufficient to prove that a solution of (1) exists and is unique; besides, it is possible to construct an iterative procedure that can guarantee a global convergence to the solution of (1) ([20]).

## 5. Convergence of the LDFI–procedure

We will now investigate the solvability of the system of nonlinear difference equations (9) by applying the LDFI–procedure when $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ and $\alpha(x, y) = 0$ in (5).

We will show that under the mild and reasonable restrictions (i)–(iv) imposed on the functions $\sigma(\boldsymbol{x}, \varphi)$ and $g(\boldsymbol{x}, \varphi)$ the problem (5)–(6) can be solved via a sequence of systems of weakly nonlinear difference equations where only $\boldsymbol{G}$ but not $\sigma$ depends on the approximate solution $\boldsymbol{u}$ of $\varphi$.

Specifically, if $\boldsymbol{u}^{(\nu)}$ is an estimate of the solution $\boldsymbol{u}^*$ of (9), we will determine a new estimate of $\boldsymbol{u}^*$ by solving the weakly nonlinear system (2)

$$\boldsymbol{F}_\nu(\boldsymbol{u}) \equiv A(\boldsymbol{u}^{(\nu)})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0}. \tag{20}$$

An approximate solution of the weakly nonlinear system (20) is computed by the simplified Newton–AM method in such a way that its solution $\boldsymbol{u}^{(\nu+1)}$ will be accepted if the residual $\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})$ satisfies the condition

$$\|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \le \varepsilon_{\nu+1}, \tag{21}$$

where $\varepsilon_{\nu+1}$ is a given tolerance such that $\varepsilon_{\nu+1} \to 0$ for $\nu \to \infty$. Here, $\|\cdot\|$ is the Euclidean norm.

If such suitable solution $\boldsymbol{u}^{(\nu+1)}$ is found, we say that the *algorithm does not break down*.

The iterate $\boldsymbol{u}^{(\nu+1)}$ is the solution of a weakly nonlinear reaction diffusion equation, whose diffusivity $\sigma$ depends on the previous iterate $\boldsymbol{u}^{(\nu)}$, with inhomogeneous term $-\boldsymbol{s} - \boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})$.

We assume that all the iterates $\boldsymbol{u}^{(\nu)}$, $\nu = 0, 1, ...$, satisfy Property A.

Thus, in particular, by inequality (12), the backward difference quotients of each grid function $u_{ij}^{(\nu)}$ are bounded. Since this bound depends on the inhomogeneous term, we have that there exist two constants $\beta > 0$ and $\beta_0 > 0$ such that

$$|\nabla_x u_{ij}^{(\nu)}| \le \beta + \varepsilon_\nu \beta_0 \quad \text{and} \quad |\nabla_y u_{ij}^{(\nu)}| \le \beta + \varepsilon_\nu \beta_0, \tag{22}$$

instead of (12), $i = 1, ..., N + 1$ and $j = 1, ..., M + 1$.

Let us prove the theorem for the convergence of the LDFI–procedure.

**Theorem 1.** Let $\boldsymbol{u}^*$ be the solution of $\boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{0}$ with $\boldsymbol{F}(\boldsymbol{u}) \equiv A(\boldsymbol{u})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s}$ arising from the discretization of the problem (5)–(6) subject to

15

the conditions (i)–(iv) and $A(\boldsymbol{u})$ being an irreducible nonsingular M–matrix and $\boldsymbol{G}(\boldsymbol{u})$ a diagonal mapping in (9).

Assume that the mapping $\boldsymbol{F}(\boldsymbol{u})$ is uniformly monotone in $\mathcal{B}$, where $\mathcal{B}$ is the set of all grid functions $\{u_{ij}\}$ that satisfy Property A (see conditions (11) and (12)).

Suppose that $\{\varepsilon_\nu\}$ is a sequence of positive numbers such that $\varepsilon_\nu \to 0$ as $\nu \to \infty$.

Let $\boldsymbol{u}^{(0)} \in \mathcal{B}$ be arbitrary and let $\boldsymbol{u}^{(\nu+1)}$ be the solution of $\boldsymbol{F}_\nu(\boldsymbol{u}) = \boldsymbol{0}$ satisfying the condition (21) with $\boldsymbol{F}_\nu(\boldsymbol{u})$ as in (20).

If all the vectors $\{\boldsymbol{u}^{(\nu)}\}$ belong to $\mathcal{B}$ and satisfy Property A with (22) instead of (12), then the sequence $\{\boldsymbol{u}^{(\nu)}\}$ converges to $\boldsymbol{u}^*$.

**Proof.** First we consider the case of $\alpha(x, y) = 0$ and $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ for the problem (5)–(6). The solution $\boldsymbol{u}^*$ in $\mathcal{B}$ of (9) satisfies the equation

$$A(\boldsymbol{u}^*)\boldsymbol{u}^* + \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{s} = \boldsymbol{0},$$

and the iterate $\boldsymbol{u}^{(\nu+1)}$ satisfies the equation

$$A(\boldsymbol{u}^{(\nu)})\boldsymbol{u}^{(\nu+1)} + \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}) - \boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}) - \boldsymbol{s} = \boldsymbol{0},$$

where the discrete $l_2(\Omega_h)$ norm of the residual $\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})$ satisfies the inequality (21).

Taking into account of the identity

$$A(\boldsymbol{u})\boldsymbol{u} - A(\boldsymbol{w})\boldsymbol{v} = A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}) + (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v},$$

for all grid functions $\boldsymbol{u}$, $\boldsymbol{v}$ and $\boldsymbol{w}$ belonging to $\mathcal{B}$, we can write

$$A(\boldsymbol{u}^*)\boldsymbol{u}^* + \boldsymbol{G}(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)})\boldsymbol{u}^{(\nu+1)} - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}) = -\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}),$$

as

$$A(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}) + (A(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)} + \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}) = -\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}).$$

Thus, we have

$$\begin{aligned} &< A(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > + \\ &+ < \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > + \\ &+ < (A(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)}, \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > = \\ &= - < \boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > . \end{aligned} \qquad (23)$$

Using (15) and assumption (ii), we can write

$$< A(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} >=$$

$$= h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \sigma_{ij}(u_{ij}^*) \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) +$$

$$+ \sum_{j=1}^{M} \sigma_{N+1j}(u_{N+1j}^*) |u_{Nj}^* - u_{Nj}^{(\nu+1)}|^2 + \sum_{i=1}^{N} \sigma_{iM+1}(u_{iM+1}^*) |u_{iM}^* - u_{iM}^{(\nu+1)}|^2$$

$$\geq \sigma_{\min} \sum_{j=1}^{M} \sum_{i=1}^{N} h^2 \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right)$$

$$+ \sigma_{\min} \left( \sum_{j=1}^{M} |u_{Nj}^* - u_{Nj}^{(\nu+1)}|^2 + \sum_{i=1}^{N} |u_{iM}^* - u_{iM}^{(\nu+1)}|^2 \right). \qquad (24)$$

Assumption (iv) on $g$ implies that, for all grid functions $\{u_{ij}\}$ and $\{v_{ij}\}$ belonging to $\mathcal{B}$, there exists a positive constant $c$ such that

$$(g_{ij}(u_{ij}) - g_{ij}(v_{ij}))(u_{ij} - v_{ij}) \geq c(u_{ij} - v_{ij})^2,$$

for all $i = 1, ..., N$ and $j = 1, ..., M$. The constant $c$ is independent of $h$. Thus for the discrete $l_2(\Omega_h)$ inner product (10) we have the inequality

$$< \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > \geq c\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2. \qquad (25)$$

Using Lemma 3 (see formula (19)) and taking into account of the assumption (iii) and the fact that, by Property A, the backward difference quotients $|\nabla_x u_{ij}^{(\nu+1)}|$ and $|\nabla_y u_{ij}^{(\nu+1)}|$ are bounded by inequalities (22), we can write

$$| < (A(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)}, \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > | \leq \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)}{2} \times$$

$$\times \sum_{j=1}^{M} \sum_{i=1}^{N} h^2 \left( \frac{|u_{ij}^* - u_{ij}^{(\nu)}|^2}{\phi} + \phi|\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) +$$

$$+ \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)}{2} \times \qquad (26)$$

$$\times \sum_{j=1}^{M} \sum_{i=1}^{N} h^2 \left( \frac{|u_{ij}^* - u_{ij}^{(\nu)}|^2}{\phi} + \phi|\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right).$$

It now follows from (23) that

$$
\begin{aligned}
< -\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} >\,\geq\, &< A(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > + \\
&+ < \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > - \\
&- |< (A(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)}, \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > |,
\end{aligned}
$$

and from (24), (25) and (26) that

$$
\sigma_{\min} \sum_{j=1}^M \sum_{i=1}^N h^2 \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) +
$$

$$
+ \sigma_{\min} \left( \sum_{j=1}^M |u_{Nj}^* - u_{Nj}^{(\nu+1)}|^2 + \sum_{i=1}^N |u_{iM}^* - u_{iM}^{(\nu+1)}|^2 \right) + c\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h -
$$

$$
- \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)}{\phi} \sum_{j=1}^M \sum_{i=1}^N h^2 |u_{ij}^* - u_{ij}^{(\nu)}|^2 -
$$

$$
- \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)\phi}{2} \sum_{j=1}^M \sum_{i=1}^N h^2 \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) \leq
$$

$$
\leq\, < -\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} >\,\leq
$$

$$
\leq \|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \, \|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h \leq \varepsilon_{\nu+1}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h,
$$

where $\phi$ is a yet an undetermined positive number.
Choosing

$$
\phi = \frac{2\sigma_{\min}}{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)},
$$

we obtain

$$
c\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2 - \frac{\Lambda^2(\beta + \varepsilon_{\nu+1}\beta_0)^2}{2\sigma_{\min}}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2 \leq \varepsilon_{\nu+1}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h. \quad (27)
$$

Since the grid function $\{u_{ij}^{(\nu+1)}\}$ belongs to $\mathcal{B}$, we may assume that

$$
\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h \leq 2\rho.
$$

Thus from (27) we have the inequality

$$
\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2 \leq \gamma\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2 + a\varepsilon_{\nu+1}, \quad (28)
$$

18

where

$$\gamma = \frac{\Lambda^2(\beta + \varepsilon_{\nu+1}\beta_0)^2}{2\sigma_{\min}c},$$

and $a = 2\rho/c$.

Now, as observed in [21], if there exists an integer $\nu_0$ such that $\gamma < 1$ for all $\nu \geq \nu_0$, we can write (28) as

$$\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu_0+\mu)}\|_h^2 \leq \gamma^\mu \|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu_0)}\|_h^2 + a\sum_{k=1}^\mu \gamma^{\mu-k}\varepsilon_{\nu_0+k},$$

$\mu = 1, 2, ...$, and since $\varepsilon_\nu \to 0$ as $\nu \to \infty$, it follows from the general Toeplitz Lemma ([26, p. 399] or [33, p. 74]) that

$$\lim_{\nu\to\infty} \|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2 = 0.$$

Therefore, the sequence $\{\boldsymbol{u}^{(\nu)}\}$ of approximate solutions converges to the solution $\boldsymbol{u}^*$ of the system (9). $\quad\sharp$

For sake of completeness, it easy to show (see [12]) that we have the convergence of $\{\boldsymbol{u}^{(\nu)}\}$ to the solution $\boldsymbol{u}^*$ of the system (9) also in the cases $\alpha(x,y) > 0$ and $\tilde{\boldsymbol{v}} \neq \boldsymbol{0}$ for the problem (5)–(6).

Indeed, since $A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{A} + \tilde{D}$, $\tilde{A}$ is the skew–symmetric part of $A(\boldsymbol{u})$; then we have

$$< A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > = < A_1(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} >,$$

and

$$< (A(\boldsymbol{u}) - A(\boldsymbol{v}))\,\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > = < (A_1(\boldsymbol{u}) - A_1(\boldsymbol{v}))\,\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >.$$

Then formula (23) becomes

$$< A_1(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > + < \tilde{D}(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > +$$
$$+ < \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > +$$
$$+ < (A_1(\boldsymbol{u}^*) - A_1(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)}, \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > =$$
$$= - < \boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} >.$$

Furthermore setting $\alpha_{\min} = \min_{(x,y)\in\Omega} \alpha(x,y)$, we have

$$< \tilde{D}(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > = h^2 \sum_{j=1}^M \sum_{i=1}^N \alpha(x_i, y_j)(u_{ij} - v_{ij})(u_{ij} - v_{ij}) \geq \alpha_{\min}\|\boldsymbol{u} - \boldsymbol{v}\|_h^2,$$

19

then, at the left hand side of inequality (27) we have to add the term $\alpha_{\min}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2$ and in (28) the parameter $\gamma$ becomes

$$\gamma = \frac{\Lambda^2(\beta + \varepsilon_{\nu+1}\beta_0)^2}{2\sigma_{\min}(\alpha_{\min} + c)}.$$

## 6. Solution of the weakly nonlinear system

In order to define the inner iterative solver for the nonlinear system (2) (or (20)), setting $\boldsymbol{w}^{(0)} = \boldsymbol{u}^{(\nu)}$, the simplified–Newton method finds the solution $\Delta\boldsymbol{w}^{(k)}$ of

$$C_\nu\Delta\boldsymbol{w} = -\boldsymbol{F}_\nu(\boldsymbol{w}^{(k)}), \tag{29}$$

for $k = 0, 1, ...$, where the matrix $C_\nu$ is the Jacobian matrix of $\boldsymbol{F}_\nu$ evaluated at the point $\boldsymbol{w}^{(0)}$, i.e., $C_\nu = F_\nu'(\boldsymbol{w}^{(0)}) = F_\nu'(\boldsymbol{u}^{(\nu)})$ and

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \Delta\boldsymbol{w}^{(k)}. \tag{30}$$

Denoting with $G'(\boldsymbol{u})$ the Jacobian matrix of $\boldsymbol{G}(\boldsymbol{u})$ that has expression

$$G'(\boldsymbol{u}) = \begin{pmatrix} \frac{\partial G_1}{\partial u_1}(u_1) & & & \\ & \frac{\partial G_2}{\partial u_2}(u_2) & & \\ & & \ddots & \\ & & & \frac{\partial G_n}{\partial u_n}(u_n) \end{pmatrix},$$

and taking into account the expression of $C_\nu = A(\boldsymbol{u}^{(\nu)}) + G'(\boldsymbol{w}^{(0)}) = A(\boldsymbol{u}^{(\nu)}) + G'(\boldsymbol{u}^{(\nu)})$ and the expression of $\boldsymbol{F}_\nu(\boldsymbol{w}^{(k)})$, formulae (29)–(30) are rewritten in such a way that the vector $\boldsymbol{w}^{(k+1)}$ is the solution of the linear system

$$C_\nu\boldsymbol{w} = G'(\boldsymbol{u}^{(\nu)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}, \tag{31}$$

for $k = 0, 1, ....$

The system (31) is solved by the block version of the Arithmetic Mean method introduced in [29].

We remind that the matrix $A(\boldsymbol{u})$ in (9) is a block tridiagonal matrix where the number of diagonal blocks is $M$.

Thus the matrix $C_\nu$ has the following form:

$$C_\nu = \begin{pmatrix} C_{11} & C_{12} & & & \\ C_{21} & C_{22} & C_{23} & & \\ & \ddots & & \ddots & \\ & & C_{MM-1} & C_{MM} \end{pmatrix}. \tag{32}$$

where $C_{ij}$ are dependent on $\boldsymbol{u}^{(\nu)}$, $i, j = 1, ..., M$.
Consider the two splittings of $C_\nu$

$$C_\nu = H_1(\boldsymbol{u}^{(\nu)}) - K_1(\boldsymbol{u}^{(\nu)}) = H_2(\boldsymbol{u}^{(\nu)}) - K_2(\boldsymbol{u}^{(\nu)}), \qquad (33)$$

where, if $M$ is even

$$H_1(\boldsymbol{u}^{(\nu)}) = \begin{pmatrix} C_{11} & C_{12} & & & & & \\ C_{21} & C_{22} & & & & & \\ & & C_{33} & C_{34} & & & \\ & & C_{43} & C_{44} & & & \\ & & & & \ddots & & \\ & & & & & C_{M-1M-1} & C_{M-1M} \\ & & & & & C_{MM-1} & C_{MM} \end{pmatrix},$$

and, consequently

$$K_1(\boldsymbol{u}^{(\nu)}) = H_1(\boldsymbol{u}^{(\nu)}) - C_\nu,$$

$$H_2(\boldsymbol{u}^{(\nu)}) = \begin{pmatrix} C_{11} & & & & & \\ & C_{22} & C_{23} & & & \\ & C_{32} & C_{33} & & & \\ & & & \ddots & & \\ & & & C_{M-2M-2} & C_{M-2M-1} & \\ & & & C_{M-1M-2} & C_{M-1M-1} & \\ & & & & & C_{MM} \end{pmatrix},$$

and

$$K_2(\boldsymbol{u}^{(\nu)}) = H_2(\boldsymbol{u}^{(\nu)}) - C_\nu.$$

If $M$ is odd, we can proceed in a similar way.
The matrices $H_1(\boldsymbol{u}^{(\nu)})$ and $H_2(\boldsymbol{u}^{(\nu)})$ are diagonally dominant and have diagonal positive entries and nonpositive off-diagonal entries; $K_1(\boldsymbol{u}^{(\nu)})$ and $K_2(\boldsymbol{u}^{(\nu)})$ are two nonnegative matrices, for all $\boldsymbol{u}^{(\nu)}$, $\nu = 0, 1, 2, ....$
Thus, the simplified Newton–AM method can be formulated as follows:

choose the initial guess $\quad \boldsymbol{w}^{(0)} = \boldsymbol{u}^{(\nu)}, \rho \geq 0;$

for $\quad k = 0, 1, ...,$ until the convergence do

$\boldsymbol{z}_k^{(0)} = \boldsymbol{w}^{(k)};$

$\quad$ for $\quad j = 1, 2, ..., j_k$ do

$\quad (H_1(\boldsymbol{u}^{(\nu)}) + \rho I)\tilde{\boldsymbol{z}}_1 = (K_1(\boldsymbol{u}^{(\nu)}) + \rho I)\boldsymbol{z}_k^{(j-1)} + (G'(\boldsymbol{u}^{(\nu)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}),$

$\quad (H_2(\boldsymbol{u}^{(\nu)}) + \rho I)\tilde{\boldsymbol{z}}_2 = (K_2(\boldsymbol{u}^{(\nu)}) + \rho I)\boldsymbol{z}_k^{(j-1)} + (G'(\boldsymbol{u}^{(\nu)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}),$

$\quad \boldsymbol{z}_k^{(j)} = \frac{1}{2}(\tilde{\boldsymbol{z}}_1 + \tilde{\boldsymbol{z}}_2);$

$\boldsymbol{w}^{(k+1)} = \boldsymbol{z}_k^{(j_k)}.$

$$\tag{34}$$

The iteration defined by the loop over $k$ will terminate when

$$\|\boldsymbol{F}_\nu(\boldsymbol{w}^{(k+1)})\| \leq \varepsilon_{\nu+1},$$

(see formula (3) or (21)). Then, $\boldsymbol{u}^{(\nu+1)} = \boldsymbol{w}^{(k+1)}$.

Here, $\{j_k\}$ denotes a sequence of positive integers. The loop over $j$ denotes the Arithmetic Mean (AM) method.

The description of the implementation and an evaluation of the effective performance of the Arithmetic Mean method on different parallel architectures are reported in the papers [29], [13], [14], [15], [16].

Let $\tilde{\boldsymbol{u}}$ be a solution of the system (2). For any vector $\boldsymbol{u}$ and $\boldsymbol{u}^{(\nu)}$ belonging to an open neighbourhood $\mathcal{K}$ of $\tilde{\boldsymbol{u}}$, we consider the following **Standard Assumptions**:

- $A(\boldsymbol{u}^{(\nu)})$ is a block tridiagonal matrix of order $n$ for any iterate $\boldsymbol{u}^{(\nu)}$.

  The diagonal blocks are square (although not necessarily all of the same order) tridiagonal submatrices, and the off–diagonal blocks are diagonal submatrices.

- The matrix $A(\boldsymbol{u}^{(\nu)})$ is irreducibly diagonally dominant and has positive diagonal entries and nonpositive off-diagonal entries for all the mesh spacings sufficiently small and for all the iterates $\boldsymbol{u}^{(\nu)} \in \mathcal{K}$.

- $\boldsymbol{G}(\boldsymbol{u})$ is a continuously differentiable diagonal mapping on $\mathbb{R}^n$ with $G'(\boldsymbol{u}) \geq 0$ for all $\boldsymbol{u} \in \mathcal{K}$.

Thus, $A(\boldsymbol{u}^{(\nu)})$ is an irreducible nonsingular M–matrix and $F'_\nu(\boldsymbol{u}) = A(\boldsymbol{u}^{(\nu)}) + G'(\boldsymbol{u})$ is also an irreducible M–matrix with $F'_\nu(\boldsymbol{u})^{-1} \leq A(\boldsymbol{u}^{(\nu)})^{-1}$ for all $\boldsymbol{u}$ and for all the iterates $\boldsymbol{u}^{(\nu)}$ belonging to $\mathcal{K}$ (see, e.g., [25, p. 109]).

We report a general result on the convergence of the simplified Newton–AM method when the Standard Assumptions are satisfied.
First we should define the matrix ($\rho \geq 0$)

$$M_\nu^{-1} = \frac{1}{2}[(H_1(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1} + (H_2(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}], \qquad (35)$$

and the iteration matrix

$$H_\nu = \frac{1}{2}[(H_1(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}(K_1(\boldsymbol{u}^{(\nu)}) + \rho I) + (H_2(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}(K_2(\boldsymbol{u}^{(\nu)}) + \rho I)], \quad (36)$$

and we observe that $H_\nu = I - M_\nu^{-1} C_\nu$.

**Theorem 2.** Suppose the system (2) $\boldsymbol{F}_\nu(\boldsymbol{u}) = \boldsymbol{0}$ has a solution $\tilde{\boldsymbol{u}}$; assume that Standard Assumptions hold for $\boldsymbol{u}$ belonging to an open neighbourhood $\mathcal{K}$ of $\tilde{\boldsymbol{u}}$ and that (33), i.e.,

$$C_\nu = H_1(\boldsymbol{u}^{(\nu)}) - K_1(\boldsymbol{u}^{(\nu)}) = H_2(\boldsymbol{u}^{(\nu)}) - K_2(\boldsymbol{u}^{(\nu)}),$$

are two splittings of the matrix $C_\nu = F'_\nu(\boldsymbol{u}^{(\nu)})$, $\boldsymbol{u}^{(\nu)} \in \mathcal{K}$, with the matrix $H_\nu$ in (36) convergent.
Then, for any $j_k \geq 1$, the solution $\tilde{\boldsymbol{u}}$ is an attraction point of the simplified Newton–AM iteration $\{\boldsymbol{w}^{(k)}\}$ defined in (34).

**Proof.** The Standard Assumptions assure that the Jacobian matrix $F'_\nu(\boldsymbol{u})$ is continuous and nonsingular and a monotone matrix in $\mathcal{K}$; in particular $C_\nu$ is a monotone matrix, $C_\nu^{-1} \geq 0$, and $H_1(\boldsymbol{u}^{(\nu)}) - K_1(\boldsymbol{u}^{(\nu)})$ and $H_2(\boldsymbol{u}^{(\nu)}) - K_2(\boldsymbol{u}^{(\nu)})$ are two weak regular splittings of $C_\nu$.
Thus, the matrices $M_\nu^{-1}$ and $H_\nu$ of (35) and (36) are nonnegative and $H_\nu$ is a convergent matrix, $\rho(H_\nu) < 1$ ([24]).
Then, from the identity

$$\left(\sum_{j=0}^{j_k-1} H_\nu^j\right)(I - H_\nu) = I - H_\nu^{j_k},$$

23

it is possible to write the simplified Newton–AM iteration (34) as

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - (\sum_{j=0}^{j_k-1} H_\nu^j) M_\nu^{-1} F_\nu(\boldsymbol{w}^{(k)}),$$

that is a *generalized linear iteration* and the proof runs as the one of 10.3.1 in [26, p. 321].    ♯

Results on the convergence and an evaluation of the effective performance of the Newton–AM method and of the simplified (or modified) Newton–AM method are reported in the papers [7], [8], [9], [10], [11].

## 7. Numerical studies

In this section we consider a numerical experimentation of the LDFI method for the solution on a rectangular domain of the model problem (5) with homogeneous Dirichlet boundary condition (6).

Different functions for the nonlinearity factors $\sigma(\boldsymbol{x}, \varphi)$ and $g(\boldsymbol{x}, \varphi)$ and for $\alpha(\boldsymbol{x})$ have been considered.

The source function $s(\boldsymbol{x})$ is chosen in order to satisfy a prespecified exact solution $\boldsymbol{u}^* = \varphi(x_j, y_j)$ of the nonlinear system (1), $i = 1, ..., N$, $j = 1, ..., M$; different choices for $\varphi(x, y)$ are examined.

In the following we list the involved functions and how they are referred. The functions $\sigma$ are dependent on $\varphi$ and are:

$$\begin{aligned}
\sigma 1 \quad &: \quad \sigma(\varphi) = 0.5 + 0.5\varphi, \\
\sigma 2 \quad &: \quad \sigma(\varphi) = 0.02 + 0.5\varphi^2, \\
\sigma 3 \quad &: \quad \sigma(\varphi) = 1/(0.02 + 0.5\varphi).
\end{aligned}$$

The functions $g$ are dependent on $\varphi$ and are [4], [5], [18], [19], [22], [23], [27], [28]:

$$\begin{aligned}
g1 \quad &: \quad g(\varphi) = 100e^{0.5\varphi}, \\
g2 \quad &: \quad g(\varphi) = -0.5e^\varphi, \\
g3 \quad &: \quad g(\varphi) = \frac{10^3\varphi}{(1 + 10\varphi)}, \\
g4 \quad &: \quad g(\varphi) = 5\varphi\log(1 + \varphi), \\
g5 \quad &: \quad g(\varphi) = 80\varphi\log(1 + \varphi).
\end{aligned}$$

We observe that

for $g1$: $g > 0$, $g' > 0$ and $g'' > 0$ for any value of $\varphi$;

for $g3$: $g \geq 0$ and $g' > 0$ when $\varphi \geq 0$;

for $g4$, $g5$: $g \geq 0$, $g' \geq 0$ and $g'' > 0$ when $\varphi \geq 0$;

and then, the functions $g1$, $g3$, $g4$ and $g5$ satisfy the Standard Assumptions on $g$ for $\varphi \geq 0$.

The chosen functions $\alpha(\boldsymbol{x})$ are the null function or:

$$\alpha 1 \ : \ \alpha(x,y) = c(x^3 + y), \qquad \text{with } c = 10, 100, 1000,$$

$$\alpha 2 \ : \ \alpha(x,y) = c\frac{10}{(10^{-3} + x + y)^2}, \qquad \text{with } c = 1, 10, 100, 1000.$$

Now we list the different functions for the exact solution.
Set

$$p(\xi) = \xi^{\hat{a}\log^2(\xi)}, \qquad q(\xi) = (2-\xi)^{\hat{a}\log^2(2-\xi)},$$

and

$$\varphi(x,y) = \begin{cases} p(x)\,p(y) & 0 < x \leq 1 & 0 < y \leq 1 \\ q(x)\,p(y) & 1 < x < 2 & 0 < y \leq 1 \\ p(x)\,q(y) & 0 < x \leq 1 & 1 < y < 2 \\ q(x)\,q(y) & 1 < x < 2 & 1 < y < 2 \\ 0 & 0 \leq x \leq 2 & y = 0, y = 2 \\ 0 & x = 0, x = 2 & 0 \leq y \leq 2 \end{cases} \tag{37}$$

then

$$\varphi 1 \ : \ \varphi(x,y) \text{ as in } (37), \qquad \hat{a} = 100,$$
$$\Omega \cup \Gamma = [0,2] \times [0,2],$$
$$\varphi 2 \ : \ \varphi(x,y) \text{ as in } (37), \qquad \hat{a} = 0.005,$$
$$\Omega \cup \Gamma = [0,2] \times [0,2].$$

Set

$$p(\xi) = \xi^{\hat{a}\log^2(\xi)}, \qquad q(\xi) = (2-\xi)^{\hat{a}\log^2(2-\xi)}, \qquad r(\xi) = -(\xi-1)^2 + 1,$$

and

$$\varphi(x,y) = \begin{cases} p(x)\,r(y) & 0 < x \leq 1 & 0 < y < 2 \\ q(x)\,r(y) & 1 < x < 2 & 0 < y < 2 \\ 0 & 0 \leq x \leq 2 & y = 0, y = 2 \\ 0 & x = 0, x = 2 & 0 \leq y \leq 2 \end{cases} \tag{38}$$

25

then

$$\varphi 3 \quad : \quad \varphi(x, y) \text{ as in (38)}, \qquad \hat{a} = 100,$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2].$$

Furthermore we have

$$\varphi 4 \quad : \quad \varphi(x, y) = \sin(\pi x) \sin(\pi y),$$
$$\Omega \cup \Gamma = [0, 1] \times [0, 1].$$

The LDFI–procedure has been implemented in a Fortran code with machine precision $2.2 \times 10^{-16}$.

In the experiments, we consider as stopping criterium for LDFI–procedure the satisfaction of both the inequalities (4)

$$\|\boldsymbol{u}^{(\nu+1)} - \boldsymbol{u}^{(\nu)}\| \le \tau_1,$$

and

$$\|\boldsymbol{F}(\boldsymbol{u}^{(\nu+1)})\| \le \tau_2,$$

with $\tau_1 = \tau_2 = 10^{-5}$.

The approximate solution computed, at each iteration of LDFI–procedure, by the simplified Newton method satisfies the stopping rule

$$\|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \le \varepsilon_{\nu+1},$$

with $\varepsilon_1 = 0.1 \, \|\boldsymbol{F}(\boldsymbol{u}^{(0)})\|$ and $\varepsilon_{\nu+1} = \min\{0.5 \, \varepsilon_\nu, \underline{\varepsilon}\}$, $\nu = 1, 2, ....$ The threshold $\underline{\varepsilon}$ is chosen $10^{-5}$, $10^{-3}$ or $10^{-2}$. In Tables 3–7, $\underline{\varepsilon}$ is chosen equal to $10^{-5}$. The starting vector of the LDFI–procedure $\boldsymbol{u}^{(0)}$ is the vector whose all components are equal to 1.

In all the experiments we have $N = M$.

In the tables, *it* indicates the number of iterations of the LDFI–procedure. The number *ktot*, the sum of the simplified Newton method's iterations, is expressed in brackets.

Here *err* denotes the computed relative error in the Euclidean norm, i.e.

$$err = \|\boldsymbol{u}^{(it)} - \boldsymbol{u}^*\| / \|\boldsymbol{u}^*\|,$$

with *res* and *res0* we indicate the residual and the initial residual in the Euclidean norm:

$$res = \|\boldsymbol{F}(\boldsymbol{u}^{(it)})\|, \qquad res0 = \|\boldsymbol{F}(\boldsymbol{u}^{(0)})\|,$$

and *diff* indicates the last difference of iterations

$$diff = \|\boldsymbol{u}^{(it)} - \boldsymbol{u}^{(it-1)}\|.$$

The term $7.60(-10)$ indicates $7.60 \times 10^{-10}$.

We indicate with $j_k$ the number of iterations of the Arithmetic Mean method for the solution of the system (31)

$$C_\nu \boldsymbol{w} = \boldsymbol{b}_\nu,$$

with $\boldsymbol{b}_\nu = G'(\boldsymbol{u}^{(\nu)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}$, that occurs at each iteration $k$ of the simplified Newton method.

At each iteration $j$, $j = 1, ..., j_k$, of the Arithmetic Mean method, $M - 1$ independent $2 \times 2$ block linear systems of order $2N$ have to be solved; the block Gaussian elimination method is used and, the parameter $\rho$ in the AM method is chosen equal to zero (see [29]).

## 8. Conclusions

In this paper we have analysed the LDFI–procedure combined with the simplified Newton–AM method for the solution of finite difference nonlinear systems.

For the convergence of the LDFI–procedure we have considered:

- a model problem where smoothness conditions on the functions involved in the equation of the model are assumed;

- the grid functions, i.e. discrete approximations of the solution of the model problem, satisfy Property A (i.e., they are uniformly bounded and have uniformly bounded backward difference quotients);

- the solution of the weakly nonlinear difference system that occurs at each iteration of the LDFI–procedure, is solved *inexactly* by a convergent iterative solver;

- the theorem of convergence of the LDFI–procedure is proved with standard techniques.

From the numerical experiments the following conclusions can be drawn:

- from Tables 1 and 2, we can observe that when the values of the function $\sigma(\varphi)$ increase ($\sigma 3$ has larger values than $\sigma 1$ and $\sigma 2$ for $\varphi \in [0,1]$), then the total number of the simplified Newton iterations increases;

- from Tables 1, 2 and 3, we observe that when the values of the function $g(\varphi)$ are rapidly increasing for $0 \leq \varphi \leq 1$ or the values of the function $\alpha(x,y)$ are large, then the diagonal of the matrix $C_\nu$ becomes more dominant; it implies a reduction of the total number of the simplified Newton iterations;

- from Tables 1 and 2, we remark that there is no appreciable reduction of the total number of the simplified Newton iterations when we solve the weakly nonlinear system with a looser accuracy than that imposed on the LDFI–procedure. Thus, in the strategy of choice of the parameters in criteria (4) and (21), these experiments suggest that the parameter $\tau_2$ must have approximately the same value of the threshold $\underline{\varepsilon}$.

- from Tables 5 and 7, we remark that the LDFI–procedure combined with the simplified Newton–AM method gives better results (in terms either of total number of simplified Newton iterations or of the number of LDFI–procedure iterations) when the coefficients of $\tilde{v}$ increase, i.e., when the *deviation from asymmetry*[2] of the matrix $C_\nu$ increases. This is a peculiar feature of the Arithmetic Mean linear solver ([29]), especially when it is implemented as inner solver in a two iteration levels procedure, such as the Newton–AM method ([9], [11]).

  In the case of the three iteration levels LDFI–procedure, the Arithmetic Mean method, as inner solver, involves a reduction of the total number of the simplified Newton iterations and an appreciable reduction of the number of the LDFI–procedure iterations when the number $j_k$ of the AM method iteration has been conveniently chosen.

## References

[1] Bai Z.Z.: *A class of two–stage iterative methods for systems of weakly nonlinear equations*, Numerical Algorithms, 14 (1997), 295–319.

---

[2]We define the deviation from asymmetry of a matrix $A$ as the difference between the Frobenius norm of the symmetric and the skew–symmetric parts of $A$ ([9]).

| $N = 256$; $\sigma(\varphi) = \sigma 2$; $\varphi(\boldsymbol{x}) = \varphi 4$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 10$; $j_k = 20$ | | | | | |
|---|---|---|---|---|---|
| | | | $\underline{\varepsilon} = 10^{-5}$ | | |
| $g(\varphi)$ | *it(ktot)* | *err* | *res* | *res0* | *diff* |
| $g1$ | 34 (174) | 7.60(-10) | 9.40(-6) | 8.08(5) | 8.10(-8) |
| $g2$ | 34 (308) | 1.80(-9) | 9.09(-6) | 8.07(5) | 1.82(-7) |
| $g3$ | 34 (266) | 1.11(-9) | 9.49(-6) | 8.09(5) | 9.85(-8) |
| $g4$ | 34 (288) | 1.63(-9) | 8.85(-6) | 8.07(5) | 2.11(-7) |
| $g5$ | 34 (167) | 7.88(-10) | 8.95(-6) | 8.08(5) | 8.87(-8) |
| | | | $\underline{\varepsilon} = 10^{-3}$ | | |
| $g(\varphi)$ | *it(ktot)* | *err* | *res* | *res0* | *diff* |
| $g1$ | 60 (175) | 7.20(-10) | 9.02(-6) | 8.08(5) | 1.46(-8) |
| $g2$ | 77 (309) | 1.93(-9) | 9.83(-6) | 8.07(5) | 2.44(-8) |
| $g3$ | 74 (268) | 1.06(-9) | 9.10(-6) | 8.09(5) | 1.43(-8) |
| $g4$ | 75 (289) | 1.70(-9) | 9.35(-6) | 8.07(5) | 2.26(-8) |
| $g5$ | 58 (167) | 8.63(-10) | 9.94(-6) | 8.08(5) | 1.84(-8) |
| | | | $\underline{\varepsilon} = 10^{-2}$ | | |
| $g(\varphi)$ | *it(ktot)* | *err* | *res* | *res0* | *diff* |
| $g1$ | 71 (175) | 7.45(-10) | 9.34(-6) | 8.08(5) | 1.52(-8) |
| $g2$ | 98 (310) | 1.82(-9) | 9.30(-6) | 8.07(5) | 2.33(-8) |
| $g3$ | 92 (268) | 1.09(-9) | 9.45(-6) | 8.09(5) | 1.51(-8) |
| $g4$ | 94 (289) | 1.77(-9) | 9.73(-6) | 8.07(5) | 2.38(-8) |
| $g5$ | 69 (168) | 7.64(-10) | 8.79(-6) | 8.08(5) | 1.64(-8) |

Table 1: Results for different functions $g(\varphi)$ and different value of $\underline{\varepsilon}$: $\sigma(\varphi) = \sigma 2$.

[2] Bai Z.Z.: *Parallel multisplitting two–stage iterative methods for large sparse systems of weakly nonlinear equations*, Numerical Algorithms, 15 (1997), 347–372.

[3] Bai Z.Z.: *Convergence analysis of the two–stage multisplitting method*, Calcolo, 36 (1999), 63–74.

[4] Brown P.N., Saad Y.: *Hybrid Krylov methods for nonlinear systems of equations*, SIAM Journal on Scientific and Statistical Computing, 11 (1990), 450–481.

[5] Buffoni G., Griffa A., Li Z., de Mottoni P.: *Spatially distributed commu-*

| $N = 256$; $\sigma(\varphi) = \sigma3$; $\varphi(\boldsymbol{x}) = \varphi4$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 10$; $j_k = 20$ | | | | | |
|---|---|---|---|---|---|
| | | | $\underline{\varepsilon} = 10^{-5}$ | | |
| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
| $g1$ | 44 (1578) | 4.77(-10) | 9.97(-6) | 7.47(7) | 1.03(-9) |
| $g2$ | 47 (3573) | 1.01(-9) | 9.94(-6) | 7.47(7) | 1.09(-9) |
| $g3$ | 45 (2660) | 8.25(-10) | 9.95(-6) | 7.47(7) | 1.12(-9) |
| $g4$ | 46 (3214) | 9.27(-10) | 9.92(-6) | 7.47(7) | 1.09(-9) |
| $g5$ | 44 (1550) | 4.56(-10) | 9.95(-6) | 7.47(7) | 1.01(-9) |
| | | | $\underline{\varepsilon} = 10^{-3}$ | | |
| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
| $g1$ | 326 (1526) | 4.68(-10) | 9.87(-6) | 7.47(7) | 9.60(-10) |
| $g2$ | 615 (3337) | 1.03(-9) | 9.97(-6) | 7.47(7) | 1.05(-9) |
| $g3$ | 499 (2530) | 8.28(-10) | 9.92(-6) | 7.47(7) | 1.06(-9) |
| $g4$ | 570 (3019) | 9.44(-10) | 9.98(-6) | 7.47(7) | 1.04(-9) |
| $g5$ | 320 (1498) | 4.94(-10) | 9.92(-6) | 7.47(7) | 9.41(-10) |
| | | | $\underline{\varepsilon} = 10^{-2}$ | | |
| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
| $g1$ | 468 (1500) | 4.69(-10) | 9.87(-6) | 7.47(7) | 9.62(-10) |
| $g2$ | 904 (3221) | 1.03(-9) | 9.96(-6) | 7.47(7) | 1.05(-9) |
| $g3$ | 728 (2466) | 8.28(-10) | 9.90(-6) | 7.47(7) | 1.05(-9) |
| $g4$ | 836 (2924) | 9.41(-10) | 9.94(-6) | 7.47(7) | 1.04(-9) |
| $g5$ | 459 (1472) | 4.52(-10) | 9.93(-6) | 7.47(7) | 9.44(-10) |

Table 2: Results for different functions $g(\varphi)$ and different value of $\underline{\varepsilon}$: $\sigma(\varphi) = \sigma3$.

nities: the resource–consumer system, Journal of Mathematical Biology, 33 (1995), 723–743.

[6] Courant R., Friedrichs K.O., Lewy H.: Über die partiellen Differenzen-gleichungen der mathematischen Physik, Mathematische Annalen, 100 (1928), 32–74.

[7] Galligani E.: A two-stage iterative method for solving a weakly nonlin-ear system, Atti del Seminario Matematico e Fisico dell'Università di Modena, L (2002), 195–215.

[8] Galligani E.: A two-stage iterative method for solving a weakly nonlinear

| $\alpha(\boldsymbol{x})$ | $c$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|---|
| 0 | | 33 (62) | 1.21(-11) | 3.26(-6) | 1.84(6) | 3.38(-9) |
| $\alpha1$ | 10 | 32 (60) | 3.72(-11) | 9.68(-6) | 1.84(6) | 9.76(-9) |
| | 100 | 31 (57) | 1.32(-11) | 3.49(-6) | 1.84(6) | 3.53(-9) |
| | 1000 | 27 (42) | 1.53(-11) | 4.98(-6) | 1.87(6) | 5.43(-9) |
| $\alpha2$ | 1 | 32 (59) | 3.70(-11) | 9.69(-6) | 1.88(6) | 9.80(-9) |
| | 10 | 29 (49) | 2.11(-11) | 5.15(-6) | 2.86(6) | 4.98(-9) |
| | 100 | 29 (32) | 1.86(-11) | 6.77(-6) | 1.94(7) | 9.06(-9) |
| | 1000 | 12 (12) | 1.81(-12) | 1.20(-6) | 1.90(8) | 3.38(-9) |

$\sigma(\varphi) = \sigma2$; $g(\varphi) = g4$; $\varphi(\boldsymbol{x}) = \varphi4$; $\tilde{v}_1 = \tilde{v}_2 = 100$; $j_k = 20$

Table 3: Results for different functions $\alpha(\boldsymbol{x})$.

| $N$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| 32 | 27 (34) | 2.26(-8) | 8.44(-6) | 1.00(4) | 6.79(-7) |
| 64 | 30 (138) | 1.17(-8) | 8.15(-6) | 5.42(4) | 3.46(-7) |
| 128 | 33 (560) | 6.95(-9) | 9.42(-6) | 2.99(5) | 9.49(-8) |
| 256 | 35 (2318) | 3.14(-9) | 8.67(-6) | 1.67(6) | 3.78(-7) |
| 512 | 38 (9510) | 1.78(-9) | 9.61(-6) | 9.43(6) | 9.33(-8) |

$\sigma(\varphi) = \sigma1$; $g(\varphi) = g4$; $\varphi(\boldsymbol{x}) = \varphi4$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 0$; $j_k = 20$

Table 4: Results for different values of $N$.

*parametrized system*, International Journal of Computer Mathematics, 79 (2002), 1211–1224.

[9] Galligani E.: *The Newton-arithmetic mean method for the solution of systems of nonlinear equations*, Applied Mathematics and Computation, 134 (2003), 9–34.

[10] Galligani E.: *The Arithmetic Mean method for solving systems of nonlinear equations in finite differences*, Applied Mathematics and Computation, 181 (2006), 579–597.

[11] Galligani E.: *On solving a special class of weakly nonlinear finite-difference systems*, International Journal of Computer Mathematics, 86 (2009), 503–522.

[12] Galligani E.: *A note on the iterative solution of nonlinear steady state*

31

| $N = 256$, $\sigma(\varphi) = \sigma1$; $g(\varphi) = g4$; $\varphi(\boldsymbol{x}) = \varphi4$; $\alpha(x,y) = 0$; $j_k = 20$ | | | | | |
|---|---|---|---|---|---|
| $\tilde{\boldsymbol{v}}$ | it(ktot) | err | res | res0 | diff |
| $(0,0)^T$ | 35 (2318) | 3.14(-9) | 8.67(-6) | 1.67(6) | 3.78(-7) |
| $(10,10)^T$ | 35 (849) | 1.18(-9) | 9.61(-6) | 1.69(6) | 1.45(-7) |
| $(50,50)^T$ | 35 (126) | 1.13(-10) | 9.88(-6) | 1.75(6) | 6.79(-9) |
| $(100,50)^T$ | 35 (74) | 2.92(-11) | 5.54(-6) | 1.79(6) | 4.65(-9) |
| $(50,100)^T$ | 34 (74) | 4.24(-11) | 8.29(-6) | 1.79(6) | 6.98(-9) |
| $(100,100)^T$ | 33 (62) | 1.21(-11) | 3.26(-6) | 1.84(6) | 3.38(-9) |
| $(150,150)^T$ | 23 (41) | 1.43(-11) | 7.55(-6) | 1.95(6) | 1.31(-8) |

Table 5: Results for different values of $\tilde{\boldsymbol{v}}$.

| $N = 256$; $\sigma(\varphi) = \sigma1$; $g(\varphi) = g4$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 100$; $j_k = 20$ | | | | | |
|---|---|---|---|---|---|
| $\varphi(\boldsymbol{x})$ | it(ktot) | err | res | res0 | diff |
| $\varphi1$ | 13 (26) | 5.32(-12) | 3.95(-7) | 5.22(5) | 4.34(-8) |
| $\varphi2$ | 26 (37) | 3.57(-12) | 7.88(-6) | 3.17(5) | 1.22(-8) |
| $\varphi3$ | 15 (28) | 2.45(-11) | 3.03(-6) | 5.23(5) | 5.83(-8) |
| $\varphi4$ | 33 (62) | 1.21(-11) | 3.26(-6) | 1.84(6) | 3.38(-9) |

Table 6: Results for different functions $\varphi(\boldsymbol{x})$.

*reaction diffusion problems*. Technical Report of Numerical Analysis, Università di Modena e Reggio Emilia, TR NA–UniMoRE–1–2010, August 2010.

[13] Galligani E., Ruggiero V.: *Computation of minimal eigenpair in the large sparse generalized eigen–problem using vector computers*, in: Parallel Computing '91 (D.J. Evans, G.R. Joubert, H. Liddell eds.), Elsevier Science Publishers B.V., North–Holland, Amsterdam, 1992, 193–201.

[14] Galligani E., Ruggiero V.: *A parallel preconditioner for block tridiagonal matrices*, in: Parallel Computing: Trends and Applications (G.R. Joubert, D. Trystram, F.J. Peters, D.J. Evans eds.), Elsevier Science Publishers B.V., North–Holland, Amsterdam, 1994, 113–120.

[15] Galligani E., Ruggiero V.: *Implementation of splitting methods for solving block tridiagonal linear systems on transputers*, in: Proceedings Euromicro Workshop on Parallel and Distributed Processing (M. Valero,

$N = 256$, $\sigma(\varphi) = \sigma 1$; $g(\varphi) = g4$; $\varphi(\boldsymbol{x}) = \varphi 4$; $\alpha(x,y) = 0$

| $\tilde{\boldsymbol{v}}$ | $j_k = 5$ | $j_k = 10$ | $j_k = 20$ | $j_k = k+1$ | $j_k = \nu+1$ |
|---|---|---|---|---|---|
| $(0,0)^T$ | 35 (9271) | 35 (4636) | 35 (2318) | 35 (1667) | 35 (2595) |
| | 3.15(-9) | 3.14(-9) | 3.14(-9) | 3.08(-9) | 3.16(-9) |
| $(10,10)^T$ | 35 (3396) | 35 (1698) | 35 (849) | 35 (1018) | 35 (1172) |
| | 1.18(-9) | 1.18(-9) | 1.18(-9) | 1.15(-9) | 1.17(-9) |
| $(50,50)^T$ | 35 (507) | 35 (254) | 35 (126) | 35 (369) | 35 (326) |
| | 1.09(-10) | 1.08(-10) | 1.13(-10) | 1.07(-10) | 9.07(-11) |
| $(100,50)^T$ | 36 (294) | 36 (148) | 35 (74) | 35 (272) | 32 (235) |
| | 4.50(-11) | 3.57(-11) | 2.92(-11) | 4.57(-11) | 2.21(-11) |
| $(50,100)^T$ | 35 (293) | 35 (147) | 34 (74) | 35 (271) | 32 (237) |
| | 4.88(-11) | 4.72(-11) | 4.24(-11) | 4.25(-11) | 3.08(-11) |
| $(100,100)^T$ | 35 (243) | 36 (122) | 33 (62) | 35 (239) | 30 (218) |
| | 3.46(-11) | 2.31(-11) | 1.21(-11) | 3.59(-11) | 1.78(-11) |
| $(150,150)^T$ | 36 (162) | 33 (82) | 23 (41) | 35 (186) | 26 (184) |
| | 1.35(-11) | 1.04(-11) | 1.43(-11) | 1.75(-11) | 1.95(-11) |

Table 7: In a column are: *it(ktot)* and *err* for different values of $j_k$.

A. Gonzalez eds.), IEEE Computer Society Press, Los Alamitos CA, 1995, 409–415.

[16] Galligani E., Ruggiero V.: *The two–stage arithmetic mean method*, Applied Mathematics and Computation, 85 (1997), 245–264.

[17] Hadamard J.: *Sur les transformations ponctuelles*, Bullettin de la Société Mathématique de France, 34 (1906), 71–84.

[18] Keller H.B., Cohen D.S.: *Some positive problems suggested by nonlinear heat generation*, Journal of Mathematics and Mechanics, 16 (1967), 1361–1376.

[19] Kernevez J.P.: Enzyme Mathematics, North–Holland, Amsterdam, 1980.

[20] Mancino O.G.: *Resolution by iteration of some nonlinear systems*, Journal of the Association for Computing Machinery, 14 (1967), 341–350.

[21] Meyer G.H.: *The numerical solution of quasilinear elliptic equations*, in: Numerical Solution of Systems of Nonlinear Algebraic Equations (G. Byrne, C.A. Hall eds.), Academic Press, New York, 1973, 27–61.

[22] Moré J.J.: *A collection of nonlinear model problems*, in: Computational Solution of Nonlinear Systems of Equations (E.L. Allgower, K. Georg eds.), Lectures in Applied Mathematics, vol. 26, American Mathematical Society, Providence RI, 1990, 723–762.

[23] Murray J.D.: Mathematical Biology, vol. I, II, Springer–Verlag, Berlin, 2003.

[24] O'Leary D.P., White R.E.: *Multi–splittings of matrices and parallel solution of linear systems*, SIAM Journal of Algebraic and Discrete Methods, 6 (1985), 630–640.

[25] Ortega J.M.: Numerical Analysis: A Second Course, Academic Press, New York, 1972. (SIAM, Philadelphia, 1990).

[26] Ortega J.M., Rheinboldt W.C.: Iterative Solution of Nonlinear Equations in Several Variables, Academic Press, New York, 1970. (SIAM, Philadelphia, 2000).

[27] Pao C.V.: *Monotone iterative methods for finite difference system of reaction diffusion equations*, Numerische Mathematik, 46 (1985), 571–586.

[28] Pao C.V.: *Nonexistence of global solutions and bifurcation analysis for a boundary value problem of parabolic type*, Proceedings of the American Mathematical Society, 65 (1977), 245–251.

[29] Ruggiero V., Galligani E.: *A parallel algorithm for solving block tridiagonal linear systems*, Computers and Mathematics with Applications, 24 (1992), 15–21.

[30] Thomas J.: Numerical Partial Differential Equations: Finite Difference Methods, Springer, New York, 1995.

[31] Varga R.S.: Matrix Iterative Analysis, Second Edition, Springer, Berlin, 2000.

[32] Wang D., Bai Z.Z., Evans D.J.: *On the monotone convergence of multisplitting method for a class of systems of weakly nonlinear equations*, International Journal of Computer Mathematics, 60 (1996), 229–242.

[33] Zygmund A.: Trigonometric Series, Second Edition, vol. I, II, Cambridge University Press, Cambridge, 1968.