# Scuola di Dottorato in
# "High Mechanics and Automotive Design & Technology"
# "Meccanica Avanzata e Tecnica del Veicolo"


# QUADERNI

## Quaderno n. 3:
## A note on the iterative solution of nonlinear steady state reaction diffusion problems

## prof. Emanuele Galligani

Meccanismo greco di Anticitera,
Museo nazionale archeologico, Atene

*All'alta fantasia qui mancò possa;*
*ma già volgeva il mio disio e 'l velle,*
*si come rota ch'igualmente è mossa,*
*l'amore che move il sole e l'altre stelle.*

Dante, Paradiso

# A note on the iterative solution of nonlinear steady state reaction diffusion problems

Emanuele Galligani

*Department of Pure and Applied Mathematics "G. Vitali"*

*Faculty of Engineering of Modena "Enzo Ferrari"*

*University of Modena–Reggio Emilia*

*Via Campi 213/b, 41100, Modena, Italy*

*III* Quaderno della Scuola

PhD School in: High Mechanics and Automotive Design and Technology

August 2010

# Riassunto

In questo quaderno si tratta la risoluzione numerica di un'equazione stazionaria di reazione e diffusione non lineare in un dominio bidimensionale.

Quando si discretizza, con il metodo alle differenze finite o agli elementi finiti, l'equazione differenziale di reazione e diffusione soggetta a condizioni al bordo di Dirichlet, si ottiene un sistema di equazioni algebriche non lineari. Nel caso della discretizzazione con differenze finite, la matrice che proviene dalla discretizzazione dei termini di diffusione e di convezione soddisfa proprietà di monotonicità.

Questo quaderno è diviso in due parti (capitoli): la prima parte riguarda la risoluzione di un'equazione di reazione e diffusione debolmente non lineare mentre la seconda parte riguarda la soluzione di un'equazione di reazione e diffusione fortemente non lineare. La soluzione di questa ultima equazione è calcolata mediante una procedura iterativa ben nota, detta LDFI, che "ritarda" la parte diffusiva dell'equazione.

Nella prima parte il sistema algebrico debolmente non lineare che proviene dalla discretizzazione dell'equazione di reazione e diffusione è risolto mediante il metodo di Newton semplificato combinato con il metodo della Media Aritmetica che è un metodo iterativo per la risoluzione di sistemi lineari sparsi di grandi dimensioni particolarmente adatto per calcolatori con architetture parallele. Questo processo di iterazione interna/esterna si può descrivere come un metodo iterativo a due stadi. Nel quaderno sono riportati risultati di convergenza globale e monotona del metodo iterativo a due stadi.

Infine, esperimenti numerici mostrano l'efficienza del metodo iterativo a due stadi, specialmente quando il termine convettivo è fortemente dominante, confermando ben noti risultati ottenuti per il caso lineare.

Nella seconda parte si analizza la procedura LDFI (Lagged Diffusivity Functional Iteration) per la risoluzione dell'equazione di reazione e diffusione nel caso fortemente non lineare.

Si considera un problema modello e si descrive una discretizzazione con un metodo alle differenze finite. Si dimostrano alcune proprietà dell'operatore alle differenze finite.

Inoltre, sono fornite condizioni sufficienti per la convergenza della procedura LDFI ad una soluzione dell'equazione di reazione e diffusione. Il sistema algebrico debolmente non lineare che si deve risolvere ad ogni passo della procedura LDFI, viene risolto con il metodo di Newton semplificato combinato con il metodo della Media Aritmetica.

Studi numerici mostrano l'efficienza della procedura LDFI combinato con il metodo di Newton semplificato-Media Aritmetica. Migliori prestazioni si ottengono quando il termine convettivo è dominante rispetto al termine diffusivo, in accordo con il caso lineare e debolmente non lineare del problema di reazione e diffusione.

**Abstract**

This report concerns with the numerical solution of nonlinear reaction diffusion equations at the steady state in a two dimensional bounded domain supplemented by suitable boundary conditions. When we use finite differences or finite element discretizations, the nonlinear diffusion equation subject to Dirichlet boundary conditions can be transcribed into a nonlinear system of algebraic equations. In the case of finite differences, the matrix that arises from the discretization of the diffusion (and/or convection) term satisfies properties of monotonicity.

This report is divided into two parts (chapters): the first part deals with the solution of a weakly nonlinear reaction diffusion equation while in the second part, the solution of a strongly nonlinear reaction diffusion equation is computed by an iterative procedure that "lags" the diffusion term. This procedure is called Lagged Diffusivity Functional Iteration (LDFI)–procedure.

In the first part the weakly nonlinear algebraic system arising from the discretization is solved by a simplified Newton method comkbined with the Arithmetic Mean method that is an iterative method, suited for parallel computers, for the solution of large sparse linear systems. This inner-outer iteration process gives a two-stage iterative method.

Results concerning the global and monotone convergence for the two-stage iterative method have been reported.

Furthermore, numerical experiments show the efficiency of the two-stage iterative method, especially for a dominant convection term, confirming the well known results for the linear case.

In the second part the LDFI-procedure for the solution of the strongly nonlinear reaction diffusion equation is analyzed.

A model problem is considered and a finite difference discretization for that model problem is described. Furthermore, in the report, 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 method is used.

Numerical studies show the efficiency for different test functions of the LDFI-procedure combined with the simplified Newton-Arithmetic Mean method. Better result are obtained for dominant convection coefficients according with the linear and the weakly nonlinear cases.

**Key Words:** Nonlinear problems, lagging diffusivity, Newton method, Arithmetic Mean method.

**AMS Classification:** 65H10, 65N06, 65N22

**C.R. Categories:** G.1.5., G.1.8.

# Chapter 1

## 1.1 Statement of the problem

A large class of relevant problems in Science and Engineering governed by reaction and diffusion processes is set up by time dependent and time independent nonlinear partial differential equations.

For example, when the reaction diffusion process reaches a steady state in a two dimensional bounded diffusion medium, we must solve an equation of the form

$$-\operatorname{div}(\sigma\nabla\varphi) + \tilde{\boldsymbol{v}} \cdot \nabla\varphi + \alpha\varphi + g(x,y,\varphi) = s(x,y) \tag{1.1}$$

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

In the equation (1.1) the convection term $\tilde{\boldsymbol{v}} \cdot \nabla\varphi$ has been taken into account; however, we will consider only problems for which convection is not dominated.

Equation (1.1) is supplemented by suitable boundary conditions on the contour $\Gamma$ of $\Omega$. The appropriate condition on the boundary depends on the physical mechanism surrounding the diffusion medium.

We assume that (1.1) has an isolated solution.

Subsequently, we report the associate expression of $g = g(\varphi)$ in (1.1) for some nonlinear reaction diffusion problems.

1. Spatially distributed communities models ([6]):

$$g(\varphi) = \frac{a\varphi^2}{(b+\varphi)} \qquad a > 0, b > 0$$

   or

$$g(\varphi) = a\varphi \log(1+\varphi) \qquad a > 0$$

2. Enzyme–substrate reaction model ([25], [30]):

$$g(\varphi) = \frac{a\varphi}{(1+b\varphi)} \qquad a, b > 0$$

3. Fischer's population growth model ([30], [34]):

$$g(\varphi) = -a\varphi(b - c\varphi) \qquad a, b, c > 0$$

   or

$$g(\varphi) = -a\varphi(\varphi - \theta)(1 - \varphi) \qquad a > 0, \quad 0 < \theta < 1$$

4. Budworm population dynamics model ([30]):

$$g(\varphi) = \frac{\varphi^2}{1+\varphi^2} - r\varphi(1 - \frac{\varphi}{q}) \qquad r, q > 0,$$

5. Molecular interaction model ([36]):

$$g(\varphi) = \varphi^2$$

6. Chemical reaction model ([1]):

$$g(\varphi) = -a(c - \varphi)e^{(-b/(1+\varphi))} \qquad a, b, c > 0$$

7. Ginzburg–Landau oscillating BZ reaction model ([22], [26])

$$g(\varphi) = -a\varphi(1 - \varphi^2) \qquad a > 0$$

and ([7])

$$g(\varphi) = -\varphi^2(1 - \varphi)$$

8. Radiation model:

$$g(\varphi) = be^{a\varphi} \qquad b \in \mathbb{R}, a > 0$$

In thermal ignition and combustion problems we have $b < 0$ and $a > 0$ (e.g. [34], [35]).

In Bratu eigenvalue problem we have $g(\varphi) = \lambda e^{\varphi}$ (e.g. [5], [29], [23]).

9. A nonlinear oscillator in an eigenvalue problem ([38])

$$g(\varphi) = -\lambda \sin \varphi$$

There is a huge amount of literature concerning the numerical solution of these initial and boundary value problems. When $\Omega$ is a two–three dimensional domain, these problems belong to the class of very large scientific computing problems for which the use of computers with parallel architecture is appropriate.

By introducing a suitable finite difference discretization scheme (see, e.g. [40, Chap. 6]), the elliptic equation (1.1) supplemented by a Dirichlet boundary condition can be transcribed into a *weakly nonlinear* system of difference equations of the form

$$\boldsymbol{F}(\boldsymbol{u}) \equiv A\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0} \tag{1.2}$$

where $\boldsymbol{u} = (u_1, u_2, ..., u_n)^T$ is a vector in $\mathbb{R}^n$, $A$ 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 $\boldsymbol{u}$.

Here, the vector $\boldsymbol{u} \in \mathbb{R}^n$ is the approximation of $\varphi(x, y)$ on a grid of points $\Omega_h$ superimposed on the domain $\Omega$. The matrix $A$, with elements $a_{ij}$, $i, j = 1, ..., n$, satisfies standard assumptions for partial difference equations of elliptic type.

**Standard assumptions:**

- $A$ is a block tridiagonal matrix of order $n$.

  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 $A$ is irreducibly diagonally dominant[1] and has positive diagonal entries and nonpositive off-diagonal entries for all the mesh spacings sufficiently small.

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

Thus, $A$ is an irreducible nonsingular M–matrix[2] and the Jacobian matrix $F'(\boldsymbol{u}) = A + G'(\boldsymbol{u})$ is also an irreducible M–matrix with $F'(\boldsymbol{u})^{-1} \leq A^{-1}$ for all $\boldsymbol{u} \in \mathbb{R}^n$ (see, e.g., [32, p. 109]).

We assume that system (1.2) has a solution $\tilde{\boldsymbol{u}}$.

---

[1]See, e.g., [40, pp. 18, 23] for definition of irreducible and irreducibly diagonally dominant matrices.
[2]A matrix $A$ of order $n$ whose elements are denoted as $a_{ij}$ is an M–matrix if $A^{-1} \geq 0$ and $a_{ij} \leq 0$ for $i \neq j$ and $i, j = 1, ..., n$ [32, p. 108].

## 1.2  A two–stage iterative method

The special form and the large dimension of the matrix $A$ in (1.2) suggest to use the simplified version of the Newton–Arithmetic Mean method for solving system (1.2) introduced in [10] (see also [9] and [13]). In order to define the iterative solver for system (1.2), setting an initial guess $\boldsymbol{w}^{(0)}$, the simplified–Newton method finds the solution $\Delta\boldsymbol{w}^{(k)}$ of

$$C\Delta\boldsymbol{w} = -\boldsymbol{F}(\boldsymbol{w}^{(k)}) \tag{1.3}$$

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

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \Delta\boldsymbol{w}^{(k)} \tag{1.4}$$

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 $C = A + G'(\boldsymbol{w}^{(0)})$ and the expression of $\boldsymbol{F}(\boldsymbol{w}^{(k)})$, formulae (1.3)–(1.4) are rewritten and then the vector $\boldsymbol{w}^{(k+1)}$ is the solution of the linear system

$$C\boldsymbol{w} = G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s} \tag{1.5}$$

for $k = 0, 1, ....$
The system (1.5) is solved by the block version of the Arithmetic Mean (AM) method introduced in ([14]).
We consider the following decomposition of the matrix[3]

$$C = \begin{pmatrix} C_{11} & C_{12} & & & \\ C_{21} & C_{22} & C_{23} & & \\ & \ddots & & \ddots & \\ & & C_{mm-1} & C_{mm} \end{pmatrix} \tag{1.6}$$

into the two splittings

$$C = H_1 - K_1 = H_2 - K_2 \tag{1.7}$$

where, if $m$ is even

$$H_1 = \begin{pmatrix} C_{11} & C_{12} & & & & & \\ C_{21} & C_{22} & & & & & \\ & & C_{33} & C_{34} & & & \\ & & C_{43} & C_{44} & & & \\ & & & & \ddots & & \\ & & & & & C_{m-1m-1} & C_{m-1m} \\ & & & & & C_{mm-1} & C_{mm} \end{pmatrix}$$

and, consequently

$$K_1 = H_1 - C$$

$$H_2 = \begin{pmatrix} C_{11} & & & & & & \\ & C_{22} & C_{23} & & & & \\ & C_{32} & C_{33} & & & & \\ & & & \ddots & & & \\ & & & & C_{m-2m-2} & C_{m-2m-1} & \\ & & & & C_{m-1m-2} & C_{m-1m-1} & \\ & & & & & & C_{mm} \end{pmatrix}$$

---

[3]If we suppose that the domain $\Omega$ is a rectangular on the $xy$ plane and we order the points of the grid $\Omega_h$ in lexicographic order, then $m$ is the number of points along the $y$ direction and each block $C_{ij}$ has order $p$, the number of points along the $x$ direction ($n = m \cdot p$).

and

$$K_2 = H_2 - C$$

If $m$ is odd, we can proceed in a similar way.

The matrices $H_1$ and $H_2$ are diagonally dominant and have diagonal positive entries and nonpositive off-diagonal entries; $K_1$ and $K_2$ are two nonnegative matrices.

Thus, the simplified Newton-Arithmetic Mean method can be formulated as follows:

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

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

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

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

$$(H_1 + \rho I)\tilde{\boldsymbol{z}}_1 = (K_1 + \rho I)\boldsymbol{z}_k^{(j-1)} + (G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s})$$

$$(H_2 + \rho I)\tilde{\boldsymbol{z}}_2 = (K_2 + \rho I)\boldsymbol{z}_k^{(j-1)} + (G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s})$$

$$\boldsymbol{z}_k^{(j)} = \tfrac{1}{2}(\tilde{\boldsymbol{z}}_1 + \tilde{\boldsymbol{z}}_2)$$ (1.8)

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

The iteration defined by the loop over $k$ will terminate when a suitable stopping criterion is satisfied, e.g., ([24, p. 73]):

$$\|\boldsymbol{F}(\boldsymbol{w}^{(k+1)})\| \leq \tau_a + \tau_r\|\boldsymbol{F}(\boldsymbol{w}^{(0)})\| \tag{1.9}$$

where $\tau_a$ and $\tau_r$ are prefixed absolute and relative error tolerances and $\|\cdot\|$ indicates a vector norm, e.g. the Euclidean norm. Then $\boldsymbol{w}^{(k+1)} \approx \tilde{\boldsymbol{u}}$.

Here, $\{j_k\}$ denotes a sequence of positive integers. The loop over $j$ denotes the Arithmetic–Mean (AM) method. This method is particularly well suited for implementation on vector–parallel computers.

An evaluation of the effective performance of the AM method on different parallel architectures is reported in the papers [14], [17], [18], [19].

If we set ($\rho \geq 0$)

$$M^{-1} = \frac{1}{2}[(H_1 + \rho I)^{-1} + (H_2 + \rho I)^{-1}] \tag{1.10}$$

$$\begin{aligned} H &= \frac{1}{2}[(H_1 + \rho I)^{-1}(K_1 + \rho I) + (H_2 + \rho I)^{-1}(K_2 + \rho I)] \\ &= I - M^{-1}C \end{aligned} \tag{1.11}$$

at each iteration $k$, the Arithmetic Mean method generates for $j = 1, 2, ..., j_k$ the vectors

$$\boldsymbol{z}_k^{(j)} = H\boldsymbol{z}_k^{(j-1)} + M^{-1}(G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s})$$

Thus,

$$\boldsymbol{w}^{(k+1)} = H^{j_k}\boldsymbol{w}^{(k)} + \left(\sum_{j=0}^{j_k-1} H^j\right) M^{-1}(G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}) \tag{1.12}$$

A helpful criterion that helps us to decide whether $H$ is convergent is provided by the following theorems ([31]).

**Proposition 1.** Let $C$ be a strictly or irreducibly diagonally dominant matrix with positive diagonal entries and nonpositive off–diagonal entries. Then, the matrix $M$ is nonsingular and the matrix $H$ is convergent.

**Proposition 2.** Let $C$ be a strictly or irreducibly diagonally dominant symmetric matrix with positive diagonal entries. Then, the matrix $M$ is nonsingular and the matrix $H$ is convergent.

We observe that, since the expression of $C$ is

$$C = A + G'(\boldsymbol{w}^{(0)})$$

setting $D = G'(\boldsymbol{w}^{(0)})$, we have the splitting for $A$

$$A = C - D$$

Then the simplified Newton–AM method can be regarded as a *two–stage iterative* method for the solution of weakly nonlinear systems ([41], [10], [2], [3], [4]).

There exist many interesting results on the convergence of the sequence $\{\boldsymbol{w}^{(k)}\}$ to a solution $\tilde{\boldsymbol{u}}$ of the system of weakly nonlinear difference equations (1.2).
First, we report the result on the convergence of Newton–iterative method when the Standard Assumptions are satisfied.

**Theorem 1.** Suppose the system (1.2) $\boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{0}$ has a solution $\tilde{\boldsymbol{u}} \in \mathbb{R}^n$; assume that Standard Assumptions hold for $\boldsymbol{u} \in \mathbb{R}^n$ (or in an open neighbourhood $\mathcal{K}$ of $\tilde{\boldsymbol{u}}$) and that (1.7), i.e.,

$$C = H_1 - K_1 = H_2 - K_2$$

are two splittings of the matrix $C = F'(\boldsymbol{w}^{(0)})$, $\boldsymbol{w}^{(0)} \in \mathbb{R}^n$ (or $\boldsymbol{w}^{(0)} \in \mathcal{K}$), with the matrix $H$ in (1.11) convergent.
Then, for any $j_k \geq 1$, the solution $\tilde{\boldsymbol{u}}$ is an attraction point of the simplified Newton–Arithmetic Mean iteration $\{\boldsymbol{w}^{(k)}\}$ defined in (1.8).

**Proof.** The Standard Assumptions assure that the Jacobian matrix $F'(\boldsymbol{u})$ is continuous and nonsingular and a monotone matrix in $\mathbb{R}^n$ (or in $\mathcal{K}$); in particular $C$ is a monotone matrix, $C^{-1} \geq 0$, and $H_1 - K_1$ and $H_2 - K_2$ are two weak regular splittings[4] of $C$. Thus, the matrices $M^{-1}$ and $H$ of (1.10) and (1.11) are nonnegative and $H$ is a convergent matrix, $\rho(H) < 1$ ([32, p. 124], [31]).
By (1.11)

$$(I - H)C^{-1} = M^{-1}$$

and from the identity

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

we obtain

$$\left(\sum_{j=0}^{j_k-1} H^j\right)M^{-1} = (I - H^{j_k})C^{-1}$$

Equation (1.12) can be written as

$$\boldsymbol{w}^{(k+1)} = H^{j_k}\boldsymbol{w}^{(k)} + (I - H^{j_k})C^{-1}(G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s})$$

Since $G'(\boldsymbol{w}^{(0)}) = C - A$ then $C^{-1}G'(\boldsymbol{w}^{(0)}) = I - C^{-1}A$, we have

$$\begin{aligned}
\boldsymbol{w}^{(k+1)} &= \boldsymbol{w}^{(k)} - (I - H^{j_k})C^{-1}(A\boldsymbol{w}^{(k)} + \boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{s}) \\
&= \boldsymbol{w}^{(k)} - \left(\sum_{j=0}^{j_k-1} H^j\right)M^{-1}(A\boldsymbol{w}^{(k)} + \boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{s})
\end{aligned}$$

thus,

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - \left(\sum_{j=0}^{j_k-1} H^j\right)M^{-1}F(\boldsymbol{w}^{(k)})$$

---

[4]$A = M - N$ is a is a weak regular splitting of the matrix $A$ if $M$ is nonsingular with $M^{-1} \geq 0$ and $M^{-1}N \geq 0$ [40, p. 95].

is a *generalized linear iteration* and the proof runs as the one of 10.3.1 in [33, p. 321].    ♯

Furthermore, important results concerning with the global and the monotone convergence of the solver (1.8) are obtained in [9] and [10].

**Theorem 2.** Let $\boldsymbol{F} : \mathbb{R}^n \to \mathbb{R}^n$ be a mapping of the form (1.2), $\boldsymbol{F}(\boldsymbol{u}) = A\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s}$, where

1. $A$ is a a monotone matrix[5], i.e., $A^{-1} \geq 0$;

2. $\boldsymbol{G} : \mathbb{R}^n \to \mathbb{R}^n$ is a positive bounded (*P–bounded*) mapping; i.e., there exists a nonnegative matrix $P$ such that

$$|\boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{G}(\boldsymbol{v})| \leq P|\boldsymbol{u} - \boldsymbol{v}|$$

   for all $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^n$ ([33, p.433])[6].

   Here $|\boldsymbol{w}|$ denotes the absolute value of the vector $\boldsymbol{w} \in \mathbb{R}^n$, i.e.,

$$|\boldsymbol{w}| = (|w_1|, |w_2|, ..., |w_n|)^T$$

3. the matrix $A^{-1}P$ has spectral radius less than one, i.e.,

$$\rho(A^{-1}P) < 1$$

4. $A = C - D$ is a regular splitting[7] of the matrix $A$;

5. $C = H_1 - K_1 = H_2 - K_2$ are two weak regular splittings of the matrix $C$;

then, for any vector $\boldsymbol{w}^{(0)}$ and for any value $j_k \geq 1$, the sequence $\{\boldsymbol{w}^{(k)}\}$ generated by the simplified Newton–AM method (1.8) converges to the solution $\tilde{\boldsymbol{u}}$ of $\boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{0}$.

**Proof.** Hypotheses 1. 2. and 3. imply that the mapping

$$\Phi(\boldsymbol{u}) = -A^{-1}(\boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s})$$

is a P–contraction ([33, p.433]), then $\boldsymbol{u} = \Phi(\boldsymbol{u})$ has a unique fixed point then $\boldsymbol{F}(\boldsymbol{u}) = 0$ has a unique solution, $\tilde{\boldsymbol{u}}$.

Since $A = C - D$ is a regular splitting we have $C^{-1} \geq 0$ and since $H_i - K_i$, $i = 1, 2$, are two weak regular splittings, thus the matrices $M^{-1}$ and $H$ of (1.10) and (1.11) are nonnegative and $H$ is convergent, i.e., $\rho(H) < 1$ ([32, p. 124], [31]). Also the matrix $D$ is nonnegative.

Writing formula (1.12) for $\tilde{\boldsymbol{u}}$ instead of $\boldsymbol{w}^{(k+1)}$ and $\boldsymbol{w}^{(k)}$ we have[8]

$$\tilde{\boldsymbol{u}} = H^{j_k}\tilde{\boldsymbol{u}} + (\sum_{j=0}^{j_k-1} H^j)M^{-1}(D\tilde{\boldsymbol{u}} - \boldsymbol{G}(\tilde{\boldsymbol{u}}) + \boldsymbol{s})$$

and subtracting this last equation from (1.12) we have

$$\boldsymbol{e}^{(k+1)} = H^{j_k}\boldsymbol{e}^{(k)} + (\sum_{j=0}^{j_k-1} H^j)M^{-1}(D\boldsymbol{e}^{(k)} - (\boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{G}(\tilde{\boldsymbol{u}})))$$

where

$$\boldsymbol{e}^{(k)} = \boldsymbol{w}^{(k)} - \tilde{\boldsymbol{u}}$$

---

[5]A matrix $A$ is a monotone matrix if $A\boldsymbol{x} \geq \boldsymbol{0}$ implies $\boldsymbol{x} \geq 0$ [21, p. 360].

[6]We remark that if $\boldsymbol{G}(\boldsymbol{u})$ is a P–bounded mapping then $\boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s}$ is also a P–bounded mapping for any vector $\boldsymbol{s}$.

[7]$A = M - N$ is a regular splitting of $A$ if $N \geq 0$ and $M$ is nonsingular and $M^{-1} \geq 0$ [32, p. 119].

[8]Formula (1.12) defines a one–step nonstationary iterative method of the kind

$$\boldsymbol{w}^{(k+1)} = \Psi_k(\boldsymbol{w}^{(k)})$$

Thus

$$|\boldsymbol{e}^{(k+1)}| \leq H^{j_k}|\boldsymbol{e}^{(k)}| + (\sum_{j=0}^{j_k-1} H^j)M^{-1}(D|\boldsymbol{e}^{(k)}| + |\boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{G}(\tilde{\boldsymbol{u}})|)$$

Since the mapping $\boldsymbol{G}$ is P–bounded we have

$$|\boldsymbol{e}^{(k+1)}| \leq S_k|\boldsymbol{e}^{(k)}|$$

where the matrix $S_k$ is defined as

$$S_k = H^{j_k} + (\sum_{j=0}^{j_k-1} H^j)M^{-1}(D + P)$$

and

$$S_k \geq 0$$

Since $A = C - D$ we have

$$S_k = H^{j_k} + (\sum_{j=0}^{j_k-1} H^j)M^{-1}C - \sum_{j=0}^{j_k-1} H^j)M^{-1}(A - P)$$

and from (1.11) we have $M^{-1}C = I - H$ and from the identity

$$(\sum_{j=0}^{j_k-1} H^j)(I - H) = I - H^{j_k}$$

we obtain

$$\begin{aligned} S_k &= H^{j_k} + (\sum_{j=0}^{j_k-1} H^j)(I - H) - (\sum_{j=0}^{j_k-1} H^j)M^{-1}(A - P) \\ &= I - (\sum_{j=0}^{j_k-1} H^j)M^{-1}(A - P) \end{aligned}$$

Since $A$ is a monotone matrix and $\rho(A^{-1}P) < 1$ we have that $A - P$ is a monotone matrix, i.e., $A - P$ is nonsingular and $(A - P)^{-1} \geq 0$. Indeed for the Neumann Lemma ([32, p. 26])

$$\begin{aligned} (A - P)^{-1} &= (A(I - A^{-1}P))^{-1} = (I - A^{-1}P)^{-1}A^{-1} \\ &= (I + A^{-1}P + (A^{-1}P)^2 + ...)A^{-1} \geq 0 \end{aligned}$$

Since the matrices $H_1 + \rho I$ and $H_2 + \rho I$ are nonsingular and $(H_1 + \rho I)^{-1} \geq 0$ and $(H_2 + \rho I)^{-1} \geq 0$ in any row of $(H_1 + \rho I)^{-1}$ and $(H_2 + \rho I)^{-1}$ there is at least one positive component; then, in any row of $M^{-1}$ there is at least one positive component. This means that $M^{-1}\boldsymbol{e} > 0$ where $\boldsymbol{e}$ is the vector whose component are all equal to 1.
If $\boldsymbol{z} = (A - P)^{-1}\boldsymbol{e}$, also $\boldsymbol{z} > 0$. Then ([2])

$$\begin{aligned} 0 &\leq S_k\boldsymbol{z} = \boldsymbol{z} - (\sum_{j=0}^{j_k-1} H^j)M^{-1}(A - P)\boldsymbol{z} \\ &= \boldsymbol{z} - (\sum_{j=0}^{j_k-1} H^j)M^{-1}\boldsymbol{e} = \boldsymbol{z} - \tilde{\boldsymbol{z}} \end{aligned}$$

with $\tilde{\boldsymbol{z}} > 0$. Clearly[9] there exists a constant $\varrho$, $0 \leq \varrho < 1$ such that $\boldsymbol{z} - \tilde{\boldsymbol{z}} \leq \varrho\boldsymbol{z}$. Thus,

$$0 \leq S_k\boldsymbol{z} \leq \varrho\boldsymbol{z} \qquad (1.13)$$

_____
[9]Since $\boldsymbol{z} - \tilde{\boldsymbol{z}} \geq 0$ with $\boldsymbol{z} > 0$ and $\tilde{\boldsymbol{z}} > 0$; let $\varrho$ be such that $0 \leq \varrho < 1$ we have for each component $i$, $i = 1, ..., n$,

$$0 \leq z_i - \tilde{z}_i \leq \varrho z_i$$

Using (1.13) repeatedly, we have

$$|\boldsymbol{e}^{(k+1)}| \leq S_k \cdot S_{k-1} \cdot ... \cdot S_0 |\boldsymbol{e}^{(0)}|$$

Let $\xi$ be such that

$$|\boldsymbol{e}^{(0)}| \leq \xi \boldsymbol{z}$$

then we obtain

$$|\boldsymbol{e}^{(k+1)}| \leq \varrho^{k+1}(\xi \boldsymbol{z})$$

that is $\lim_{k \to \infty} \boldsymbol{w}^{(k)} = \tilde{\boldsymbol{u}}$.        ♯

**Theorem 3.** Let $\boldsymbol{F} : \mathbb{R}^n \to \mathbb{R}^n$ be a mapping of the form (1.2), $\boldsymbol{F}(\boldsymbol{u}) = A\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s}$; assume that a solution $\tilde{\boldsymbol{u}}$ of the system of equations $\boldsymbol{F}(\boldsymbol{u}) = 0$ exists in $\mathbb{R}^n$ and that

1. the Jacobian matrix $F'(\boldsymbol{u}) = A + G'(\boldsymbol{u})$ is a monotone matrix for all $\boldsymbol{u} \in \mathbb{R}^n$, i.e., $F'(\boldsymbol{u})^{-1} \geq 0$;

2. $\boldsymbol{G}(\boldsymbol{u})$ is a continuously differentiable diagonal mapping on $\mathbb{R}^n$ and, in addition, $G'(\boldsymbol{u}) \geq 0$ for all $\boldsymbol{u} \in \mathbb{R}^n$ and $G'(\boldsymbol{u})$ satisfies the isotonicity condition $G'(\boldsymbol{u}) \leq G'(\boldsymbol{v})$ whenever $\boldsymbol{u} \leq \boldsymbol{v}$;

3. let $C$ be the matrix $C = F'(\boldsymbol{w}^{(0)})$, where $\boldsymbol{w}^{(0)}$ is a point of $\mathbb{R}^n$ satisfying the condition $\boldsymbol{F}(\boldsymbol{w}^{(0)}) \geq 0$;

4. let

$$C = H_1 - K_1 = H_2 - K_2$$

be two weak reagular splittings of the matrix $C$;

then[10] the sequence $\{\boldsymbol{w}^{(k)}\}$ generated by the simplified Newton–AM method (1.8), starting from $\boldsymbol{w}^{(0)}$, with the number of iterations $j_k = J$ fixed at each stage $k$, satisfies

$$\tilde{\boldsymbol{u}} \leq \boldsymbol{w}^{(k+1)} \leq \boldsymbol{w}^{(k)} \qquad k = 0, 1, 2, ... \tag{1.14}$$

**Proof.** By hypothesis $C$ is a monotone matrix, $C^{-1} \geq 0$, and $H_1 - K_1$ and $H_2 - K_2$ are two weak regular splittings of $C$. Thus, the matrices $M^{-1}$ and $H$ of (1.10) and (1.11) are nonnegative and $H$ is a convergent matrix, $\rho(H) < 1$ ([32, p. 124], [31]). Also the matrix $G'(\boldsymbol{w}^{(0)})$ is nonnegative.
By (1.11)

$$(I - H)C^{-1} = M^{-1}$$

and from the identity

$$(\sum_{j=0}^{J-1} H^j)(I - H) = I - H^J$$

we obtain

$$(\sum_{j=0}^{J-1} H^j)M^{-1} = (I - H^J)C^{-1} \tag{1.15}$$

Equation (1.12) can be written as

$$\boldsymbol{w}^{(k+1)} = H^J \boldsymbol{w}^{(k)} + (I - H^J)C^{-1}(G'(\boldsymbol{w}^{(0)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s}) \tag{1.16}$$

that implies

$$(1 - \varrho)z_i \leq \tilde{z}_i$$

Thus

$$(1 - \varrho)z_{\max} \leq \tilde{z}_{\min}$$

we obtain

$$\varrho = \frac{z_{\max} - \tilde{z}_{\min}}{z_{\max}}$$

[10]We remark that if $A$ is an M–matrix and $\boldsymbol{G}(\boldsymbol{u})$ is continuously differentiable diagonal mapping on $\mathbb{R}^n$ with $G'(\boldsymbol{u}) \geq 0$ for all $\boldsymbol{u} \in \mathbb{R}^n$ then $F'(\boldsymbol{u})$ is an M–matrix for any $\boldsymbol{u} \in \mathbb{R}^n$ and the solution of $\boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{0}$ exists and is unique [33, p. 141].

Since $G'(\boldsymbol{w}^{(0)}) = C - A$ then $C^{-1}G'(\boldsymbol{w}^{(0)}) = I - C^{-1}A$, we have

$$
\begin{aligned}
\boldsymbol{w}^{(k+1)} &= \boldsymbol{w}^{(k)} - (I - H^J)C^{-1}(A\boldsymbol{w}^{(k)} + \boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{s}) \\
&= \boldsymbol{w}^{(k)} - (\sum_{j=0}^{J-1} H^j)M^{-1}(A\boldsymbol{w}^{(k)} + \boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{s})
\end{aligned}
\tag{1.17}
$$

Since $\tilde{\boldsymbol{u}}$ is a solution of (1.2) and $A = C - G'(\boldsymbol{w}^{(0)})$ we have

$$
A\tilde{\boldsymbol{u}} + \boldsymbol{G}(\tilde{\boldsymbol{u}}) - \boldsymbol{s} = \boldsymbol{0} \quad \Longrightarrow \quad (C - G'(\boldsymbol{w}^{(0)}))\tilde{\boldsymbol{u}} + \boldsymbol{G}(\tilde{\boldsymbol{u}}) - \boldsymbol{s} = \boldsymbol{0}
$$

Thus,

$$
\tilde{\boldsymbol{u}} = C^{-1}(G'(\boldsymbol{w}^{(0)})\tilde{\boldsymbol{u}} - \boldsymbol{G}(\tilde{\boldsymbol{u}}) + \boldsymbol{s})
$$

and writing (1.16) for $\tilde{\boldsymbol{u}}$ instead of $\boldsymbol{w}^{(k)}$ and $\boldsymbol{w}^{(k+1)}$ we have

$$
\tilde{\boldsymbol{u}} = C^{-1}(G'(\boldsymbol{w}^{(0)})\tilde{\boldsymbol{u}} - \boldsymbol{G}(\tilde{\boldsymbol{u}}) + \boldsymbol{s}) + H^J(\tilde{\boldsymbol{u}} - C^{-1}(G'(\boldsymbol{w}^{(0)})\tilde{\boldsymbol{u}} - \boldsymbol{G}(\tilde{\boldsymbol{u}}) + \boldsymbol{s}))
$$

Subtracting this last equation from (1.16) we obtain

$$
\boldsymbol{w}^{(k+1)} - \tilde{\boldsymbol{u}} = H^J(\boldsymbol{w}^{(k)} - \tilde{\boldsymbol{u}}) + (I - H^J)C^{-1}(G'(\boldsymbol{w}^{(0)})(\boldsymbol{w}^{(k)} - \tilde{\boldsymbol{u}}) - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{G}(\tilde{\boldsymbol{u}}))
$$

or for (1.17)

$$
\boldsymbol{e}^{(k+1)} = H^J\boldsymbol{e}^{(k)} + (\sum_{j=0}^{J-1} H^j)M_\nu^{-1}(G'(\boldsymbol{w}^{(0)})\boldsymbol{e}^{(k)} - (\boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{G}(\tilde{\boldsymbol{u}})))
\tag{1.18}
$$

where $\boldsymbol{e}^{(k)} = \boldsymbol{w}^{(k)} - \tilde{\boldsymbol{u}}$.
Using the Mean Value Theorem (e.g., [33, p. 68])

$$
\boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{G}(\tilde{\boldsymbol{u}}) = G'(\boldsymbol{\xi}^{(k)})(\boldsymbol{w}^{(k)} - \tilde{\boldsymbol{u}})
\tag{1.19}
$$

where $\boldsymbol{\xi}^{(k)} \in (\tilde{\boldsymbol{u}}, \boldsymbol{w}^{(k)})$ and $G'(\boldsymbol{\xi}^{(k)})$ is a diagonal matrix. By hypothesis $G'(\boldsymbol{\xi}^{(k)}) \geq 0$.

Since the hypotheses on the matrices $H_i$ and $K_i$, $i = 1, 2$, the matrix

$$
(\sum_{j=0}^{J-1} H^j)M^{-1}
$$

is nonnegative. Furthermore, by the hypothesis $\boldsymbol{F}(\boldsymbol{w}^{(0)}) \geq 0$, thus, for $k = 0$ formula (1.17) gives

$$
\boldsymbol{w}^{(1)} - \boldsymbol{w}^{(0)} \leq \boldsymbol{0}
$$

Also

$$
\tilde{\boldsymbol{u}} - \boldsymbol{w}^{(0)} \leq \boldsymbol{0}
$$

Indeed,

$$
\boldsymbol{0} \leq \boldsymbol{F}(\boldsymbol{w}^{(0)}) = \boldsymbol{F}(\boldsymbol{w}^{(0)}) - \boldsymbol{F}(\tilde{\boldsymbol{u}}) = F'(\boldsymbol{\xi}^{(0)})(\boldsymbol{w}^{(0)} - \tilde{\boldsymbol{u}})
$$

By hypothesis $F'(\boldsymbol{\xi}^{(0)})$ is a monotone matrix, thus $\boldsymbol{w}^{(0)} - \tilde{\boldsymbol{u}} \geq \boldsymbol{0}$.
Using (1.19) the term $G'(\boldsymbol{w}^{(0)})\boldsymbol{e}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) - \boldsymbol{G}(\tilde{\boldsymbol{u}})$ can be rewritten for $k = 0$ as

$$
G'(\boldsymbol{w}^{(0)})\boldsymbol{e}^{(0)} - G'(\boldsymbol{\xi}^{(0)})(\boldsymbol{w}^{(0)} - \tilde{\boldsymbol{u}}) \quad \Longrightarrow \quad (G'(\boldsymbol{w}^{(0)}) - G'(\boldsymbol{\xi}^{(0)}))(\boldsymbol{w}^{(0)} - \tilde{\boldsymbol{u}})
$$

The assumption on isotonicity condition of $G'$ implies $G'(\boldsymbol{w}^{(0)}) - G'(\boldsymbol{\xi}^{(0)}) \geq 0$, since $\boldsymbol{\xi}^{(0)} \in (\tilde{\boldsymbol{u}}, \boldsymbol{w}^{(0)})$ and $\tilde{\boldsymbol{u}} \leq \boldsymbol{w}^{(0)}$; then

$$
(G'(\boldsymbol{w}^{(0)}) - G'(\boldsymbol{\xi}^{(0)}))(\boldsymbol{w}^{(0)} - \tilde{\boldsymbol{u}}) \geq \boldsymbol{0}
$$

Thus, for $k = 0$, formula (1.18) gives

$$
\boldsymbol{w}^{(1)} - \tilde{\boldsymbol{u}} \geq \boldsymbol{0} \quad \Longrightarrow \quad \tilde{\boldsymbol{u}} \leq \boldsymbol{w}^{(1)} \leq \boldsymbol{w}^{(0)}
$$

Since in (1.16) $j$ is fixed for any $k$, by using (1.15), we can rewrite formula (1.16) for $\boldsymbol{w}^{(k)}$ and $\boldsymbol{w}^{(k-1)}$

$$\boldsymbol{w}^{(k)} = H^J \boldsymbol{w}^{(k-1)} + (\sum_{j=0}^{J-1} H^j) M^{-1} (G'(\boldsymbol{w}^{(0)}) \boldsymbol{w}^{(k-1)} - \boldsymbol{G}(\boldsymbol{w}^{(k-1)}) + \boldsymbol{s}) \tag{1.20}$$

Subtracting (1.16) from (1.20) we have

$$\boldsymbol{w}^{(k)} - \boldsymbol{w}^{(k+1)} = H^J (\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^{(k)}) + (\sum_{j=0}^{J-1} H^j) M_\nu^{-1} G'(\boldsymbol{w}^{(0)}) (\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^{(k)}) -$$

$$- (\sum_{j=0}^{J-1} H^j) M^{-1} (\boldsymbol{G}(\boldsymbol{w}^{(k-1)}) - \boldsymbol{G}(\boldsymbol{w}^{(k)})) \tag{1.21}$$

Since $\boldsymbol{w}^{(0)} \geq \boldsymbol{w}^{(1)}$ we can write

$$\boldsymbol{G}(\boldsymbol{w}^{(0)}) - \boldsymbol{G}(\boldsymbol{w}^{(1)}) = G'(\bar{\boldsymbol{\xi}}^{(0)})(\boldsymbol{w}^{(0)} - \boldsymbol{w}^{(1)})$$

where $\bar{\boldsymbol{\xi}}^{(0)} \in (\boldsymbol{w}^{(1)}, \boldsymbol{w}^{(0)})$.
Thus, formula (1.21) for $k = 1$ becomes

$$\boldsymbol{w}^{(1)} - \boldsymbol{w}^{(2)} = H^J (\boldsymbol{w}^{(0)} - \boldsymbol{w}^{(1)}) + (\sum_{j=0}^{J-1} H^j) M^{-1} \cdot$$

$$\cdot (G'(\boldsymbol{w}^{(0)}) - G'(\bar{\boldsymbol{\xi}}^{(0)}))(\boldsymbol{w}^{(0)} - \boldsymbol{w}^{(1)})$$

Using the isotonicity condition on $G'$, it follows that $\boldsymbol{w}^{(1)} - \boldsymbol{w}^{(2)} \geq \boldsymbol{0}$.
Formula (1.18) with (1.19) for $k = 1$ gives $\boldsymbol{w}^{(2)} - \tilde{\boldsymbol{u}} \geq \boldsymbol{0}$, since $\boldsymbol{\xi}^{(1)} \in (\tilde{\boldsymbol{u}}, \boldsymbol{w}^{(1)})$ and $G'(\boldsymbol{w}^{(0)}) \geq G'(\boldsymbol{\xi}^{(1)})$.
We have thus shown that

$$\tilde{\boldsymbol{u}} \leq \boldsymbol{w}^{(2)} \leq \boldsymbol{w}^{(1)} \leq \boldsymbol{w}^{(0)}$$

By induction, using formula (1.21) with

$$\boldsymbol{G}(\boldsymbol{w}^{(k-1)}) - \boldsymbol{G}(\boldsymbol{w}^{(k)}) = G'(\bar{\boldsymbol{\xi}}^{(k-1)})(\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^{(k)})$$

where $\bar{\boldsymbol{\xi}}^{(k-1)} \in (\boldsymbol{w}^{(k)}, \boldsymbol{w}^{(k-1)})$, $G'(\bar{\boldsymbol{\xi}}^{(k-1)}) \leq G'(\boldsymbol{w}^{(0)})$, and formula (1.18) with (1.19), we may complete the proof of the theorem. $\quad \sharp$

We observe that Theorem 3 holds also in the case of $\boldsymbol{u} \in \mathcal{K}$, where $\mathcal{K}$ is a compact set of $\mathbb{R}^n$ neighbourhood of $\tilde{\boldsymbol{u}}$ and with $\boldsymbol{w}^{(0)} \in \mathcal{K}$.
Furthermore, we note that the hypothesis 2. on the mapping $\boldsymbol{G}(\boldsymbol{u})$ can be replaced by a condition of one–sided positive boundedness on $\boldsymbol{G}$.

**Definition.** A mapping $\boldsymbol{G} : \mathbb{R}^n \to \mathbb{R}^n$ is one–sided positive bounded (*one–sided P–bounded*) if there exists a nonnegative matrix $P$ such that

$$\boldsymbol{0} \leq \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{G}(\boldsymbol{v}) \leq P(\boldsymbol{u} - \boldsymbol{v})$$

for all $\boldsymbol{u} \geq \boldsymbol{v}$ ([9]).

## 1.3  Numerical experiments

In this section we consider a numerical experimentation of the simplified Newton Arithmetic Mean method for the solution on a square domain of the problem (1.1) with homogeneous Dirichlet boundary conditions and with nonhomogeneous Dirichlet boundary conditions

$$\varphi(\boldsymbol{x}) = U_0(\boldsymbol{x}) \qquad \boldsymbol{x} \in \Gamma \tag{1.22}$$

In the case of nonhomogeneous boundary conditions (1.22) the system (1.2) becomes

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

where the vector $\boldsymbol{b}$ depends on the values of the solution at points of the discretization of the boundary $\Gamma$.

Different functions for the nonlinearity factor $g(\boldsymbol{x}, \varphi)$ have been considered and are refereed as

$$
\begin{aligned}
g1_1 \quad &: \quad g(\varphi) = e^{\varphi} \\
g1_2 \quad &: \quad g(\varphi) = 100e^{0.5\varphi} \\
g1_3 \quad &: \quad g(\varphi) = -0.5e^{\varphi} \\
g2_1 \quad &: \quad g(\varphi) = \frac{\varphi}{(1 + \varphi)} \\
g2_2 \quad &: \quad g(\varphi) = \frac{1000\varphi}{(1 + 10\varphi)} \\
g3 \quad &: \quad g(\varphi) = -(0.4 - \varphi)e^{(-15/(1+\varphi))} \\
g4 \quad &: \quad g(\varphi) = \frac{0.02\varphi^2}{(3 + \varphi)} \\
g5_1 \quad &: \quad g(\varphi) = 0.005\varphi \log(1 + \varphi) \\
g5_2 \quad &: \quad g(\varphi) = 5\varphi \log(1 + \varphi) \\
g5_3 \quad &: \quad g(\varphi) = 500\varphi \log(1 + \varphi) \\
g5_4 \quad &: \quad g(\varphi) = 80\varphi \log(1 + \varphi) \\
g6 \quad &: \quad g(\varphi) = -\varphi(2 - \varphi)
\end{aligned}
$$

We observe that

for $g1_1$, $g1_2$: $g > 0$, $g' > 0$ and $g'' > 0$ for any value of $\varphi$;

for $g2_k$ $(k = 1, 2)$: $g \geq 0$ and $g' > 0$ when $\varphi \geq 0$;

for $g3$: $g \geq 0$, $g' \geq 0$ and $g'' \geq 0$ when $\varphi \geq 0.4$;

for $g4$ and $g5_k$ $(k = 1, 4)$: $g \geq 0$, $g' \geq 0$ and $g'' > 0$ when $\varphi \geq 0$;

and then, only for the functions $g1_1$, $g1_2$, $g4$ and $g5_k$, $(k = 1, ..., 4)$, the conditions on $g$ of the Theorem 3 are satisfied for $\varphi \geq 0$.

The source function $s(\boldsymbol{x})$ is chosen in order to satisfy the system (1.2) with a prespecified exact solution $\tilde{\boldsymbol{u}}$. Here the components of $\tilde{\boldsymbol{u}}$ are the values of the function $\varphi(\boldsymbol{x})$ at the grid points $\Omega_h$; different choices for $\varphi(\boldsymbol{x})$ are examined.

We list different functions for the exact solution $\varphi(x, y)$ of (1.1):

$$
\begin{aligned}
\varphi1 \quad &: \quad \varphi(x, y) = (x - a)(y - a)(b - x)(b - y)(2x^2 + y) \qquad a = 0 \quad b = 1 \\
&\qquad \Omega \cup \Gamma = [0, 1] \times [0, 1]
\end{aligned}
$$

Set

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

and

$$
\varphi(x, y) = \begin{cases}
p(x) \cdot p(y) & 0 < x \leq 1 & 0 < y \leq 1 \\
q(x) \cdot p(y) & 1 < x < 2 & 0 < y \leq 1 \\
p(x) \cdot q(y) & 0 < x \leq 1 & 1 < y < 2 \\
q(x) \cdot q(y) & 1 < x < 2 & 1 < y < 2 \\
0 & 0 \leq x \leq 2 & y = 0, y = 2 \\
0 & x = 0, x = 2 & 0 \leq y \leq 2
\end{cases} \tag{1.23}
$$

then

$$\begin{aligned}
\varphi2_1 \quad &: \quad \varphi(x,y) \text{ as in } (1.23); \qquad \hat{\alpha} = 100 \\
&\qquad \Omega \cup \Gamma = [0,2] \times [0,2] \\
\varphi2_2 \quad &: \quad \varphi(x,y) \text{ as in } (1.23); \qquad \hat{\alpha} = 0.05 \\
&\qquad \Omega \cup \Gamma = [0,2] \times [0,2] \\
\varphi2_3 \quad &: \quad \varphi(x,y) \text{ as in } (1.23); \qquad \hat{\alpha} = 0.005 \\
&\qquad \Omega \cup \Gamma = [0,2] \times [0,2]
\end{aligned}$$

Set

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

and

$$\varphi(x,y) = \begin{cases} p(x) \cdot r(y) & 0 < x \le 1 & 0 < y < 2 \\ q(x) \cdot r(y) & 1 < x < 2 & 0 < y < 2 \\ 0 & 0 \le x \le 2 & y = 0, y = 2 \\ 0 & x = 0, x = 2 & 0 \le y \le 2 \end{cases} \qquad (1.24)$$

then

$$\begin{aligned}
\varphi3_1 \quad &: \quad \varphi(x,y) \text{ as in } (1.24); \qquad \hat{\alpha} = 100 \\
&\qquad \Omega \cup \Gamma = [0,2] \times [0,2] \\
\varphi3_2 \quad &: \quad \varphi(x,y) \text{ as in } (1.24); \qquad \hat{\alpha} = 0.05 \\
&\qquad \Omega \cup \Gamma = [0,2] \times [0,2]
\end{aligned}$$

Furthermore we have

$$\begin{aligned}
\varphi4 \quad &: \quad \varphi(x,y) = (1 + x - y)^3 \\
&\qquad \Omega \cup \Gamma = [0,1] \times [0,1]
\end{aligned}$$

We notice that the functions $\varphi1$, $\varphi2_1$, $\varphi2_2$, $\varphi2_3$, $\varphi3_1$ and $\varphi3_1$ are equal to zero at the points of the boundary $\Gamma$.

The functions $\sigma$ and $\alpha$ of (1.1) are taken as $\sigma(x,y) = 1$ and $\alpha(x,y) = 0$ or $\alpha(x,y) = (x^3 + y)$.
Different values of the terms multiplying first–order derivatives, $\tilde{v}_1$ and $\tilde{v}_2$, are considered.

The simplified Newton–AM method has been implemented in a Fortran code with machine precision $2.2 \times 10^{-16}$. The absolute and relative tolerance for the simplified Newton method stopping criterium (1.9) are taken equal to $10^{-10}$. The number of inner iterations $j_k$ of the Arithmetic Mean method (1.8) is prefixed equal to 20.
In the experiments the domain $\Omega$ is a square domain and we superimpose a grid that has the number of points along the $x$–axis equal to the one along the $y$–axis; we denote $m$ this number. Here $n = m \cdot m$ and the mesh width is $h = 1/(m + 1)$.
When $|\tilde{v}_1|h < 2$ and $|\tilde{v}_2|h < 2$, the matrix $A$ is a irreducibly diagonally dominant M–matrix and satisfies Standard Assumptions.
The starting vector of the method $\boldsymbol{w}^{(0)}$ is the vector whose all components are equal to 1.
In the tables, *it* indicates the number of iterations of the simplified Newton–AM method needed to satisfy the stopping rule (1.9), *err* denotes the computed relative error in the Euclidean norm, i.e.

$$err = \|\boldsymbol{w}^{(it)} - \tilde{\boldsymbol{u}}\| / \|\tilde{\boldsymbol{u}}\|$$

and *res* and *res0* are the residual and the initial residual in the Euclidean norm:

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

The term $8.35(-9)$ indicates $8.35 \cdot 10^{-9}$.

$\varphi(\boldsymbol{x}) = \varphi 1,\ g(\varphi) = g1_1,\ \alpha(x,y) = 0$

| | | $\tilde{v}_1 = \tilde{v}_2 = 0$ | | | $\tilde{v}_1 = \tilde{v}_2 = 100$ | |
|---|---|---|---|---|---|---|
| *m* | *it* | *err* | *res* | *it* | *err* | *res* |
| 32 | 54 | 8.35(-9) | 9.52(-7) | 5 | 4.77(-12) | 1.58(-8) |
| 64 | 193 | 2.85(-8) | 6.84(-6) | 6 | 2.51(-11) | 5.18(-7) |
| 128 | 714 | 8.15(-8) | 3.99(-5) | 15 | 1.60(-10) | 1.40(-5) |
| 256 | 2665 | 2.28(-7) | 2.26(-4) | 55 | 1.15(-9) | 1.98(-4) |

$\varphi(\boldsymbol{x}) = \varphi 1,\ g(\varphi) = g1_1,\ \alpha(x,y) = (x^3 + y)$

| | | | | | | |
|---|---|---|---|---|---|---|
| 128 | 693 | 7.87(-8) | 3.98(-5) | 15 | 1.58(-10) | 1.38(-5) |

Table 1.1: Results for different values of *m*.

$m = 128,\ g(\varphi) = g1_1,\ \alpha(x,y) = 0$

| | | $\tilde{v}_1 = \tilde{v}_2 = 0$ | | | $\tilde{v}_1 = \tilde{v}_2 = 100$ | |
|---|---|---|---|---|---|---|
| $\varphi(\boldsymbol{x})$ | *it* | *err* | *res* | *it* | *err* | *res* |
| $\varphi 1$ | 713 | 3.64(-7) | 3.76(-5) | 15 | 7.57(-10) | 1.39(-5) |
| $\varphi 2_1$ | 612 | 8.50(-8) | 9.38(-6) | 9 | 9.15(-12) | 2.81(-7) |
| $\varphi 2_2$ | 426 | 1.26(-8) | 1.00(-5) | 8 | 1.96(-11) | 4.38(-6) |
| $\varphi 2_3$ | 391 | 4.24(-9) | 3.99(-6) | 8 | 2.17(-13) | 8.48(-8) |
| $\varphi 3_1$ | 583 | 4.10(-8) | 9.21(-6) | 9 | 3.46(-12) | 2.09(-7) |
| $\varphi 3_2$ | 475 | 1.51(-8) | 9.41(-6) | 9 | 1.23(-13) | 2.02(-8) |
| $\varphi 4$ | 714 | 8.15(-8) | 3.99(-5) | 15 | 1.60(-10) | 1.40(-5) |

Table 1.2: Results for different values of $\varphi(\boldsymbol{x})$.

The experiments highlight the property of the method that it gives better results when applied to non symmetric problems. This property also happens for the Arithmetic Mean method to linear problems. Furthermore, 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 $\alpha(x,y)$ increase, then the diagonal of the matrix $C$ in (1.3) becomes more dominant; it implies a reduction of the number of the simplified Newton iterations.

$m = 128,\ \varphi(\boldsymbol{x}) = \varphi 2_3,\ \alpha(x,y) = 0$

| | $\tilde{v}_1 = \tilde{v}_2 = 0$ | | | $\tilde{v}_1 = \tilde{v}_2 = 100$ | | |
|---|---|---|---|---|---|---|
| $g(\varphi)$ | $it$ | $err$ | $res$ | $it$ | $err$ | $res$ |
| $g1_1$ | 391 | 4.24(-9) | 3.99(-6) | 8 | 2.17(-13) | 8.48(-8) |
| $g1_2$ | 41 | 3.09(-10) | 3.32(-6) | 9 | 2.95(-13) | 1.21(-7) |
| $g1_3$ | 789 | 9.28(-9) | 4.10(-6) | 8 | 2.96(-14) | 1.21(-8) |
| $g2_1$ | 560 | 6.35(-9) | 4.06(-6) | 8 | 1.05(-14) | 7.81(-10) |
| $g2_2$ | 234 | 2.48(-9) | 4.06(-6) | 9 | 8.64(-14) | 4.86(-8) |
| $g3$ | 586 | 6.69(-9) | 4.07(-6) | 8 | 5.74(-15) | 1.48(-10) |
| $g4$ | 585 | 6.72(-9) | 4.10(-6) | 8 | 1.10(-14) | 2.03(-10) |
| $g5_1$ | 585 | 6.78(-9) | 4.13(-6) | 8 | 1.82(-14) | 3.03(-10) |
| $g5_2$ | 281 | 3.04(-9) | 4.08(-6) | 8 | 6.91(-13) | 2.91(-7) |
| $g5_3$ | 9 | 7.24(-12) | 6.59(-7) | 8 | 1.09(-11) | 3.71(-6) |
| $g5_4$ | 36 | 2.82(-10) | 3.49(-6) | 9 | 1.80(-12) | 8.01(-7) |
| $g6$ | 587 | 6.69(-9) | 4.06(-6) | 8 | 1.48(-13) | 6.21(-8) |

Table 1.3: Results for different values of $g(\varphi)$.

$m = 128,\ \varphi(\boldsymbol{x}) = \varphi 1,\ g(\varphi) = g1_1,\ \alpha(x,y) = 0$

| $\tilde{v}_1$ | $\tilde{v}_2$ | $it$ | $err$ | $res$ | $res0$ |
|---|---|---|---|---|---|
| 0 | 0 | 714 | 8.15(-8) | 3.99(-5) | 4.05(5) |
| 10 | 10 | 254 | 2.41(-8) | 4.02(-5) | 4.06(5) |
| 50 | 100 | 18 | 3.11(-10) | 1.85(-5) | 4.08(5) |
| 100 | 50 | 19 | 2.43(-10) | 1.40(-5) | 4.39(5) |
| 100 | 100 | 15 | 1.60(-10) | 1.40(-5) | 4.34(5) |
| 257 | 257 | 7 | 8.73(-11) | 1.70(-5) | 5.68(5) |
| 300 | 300 | 7 | 7.13(-13) | 3.62(-7) | 6.17(5) |
| 350 | 350 | 7 | 1.41(-12) | 1.34(-6) | 6.77(5) |
| 400 | 400 | 7 | 5.17(-11) | 6.28(-5) | 7.41(5) |

Table 1.4: Results for different values of $\tilde{v}_1$ and $\tilde{v}_2$.



$\varphi(x,y) = \varphi 1$      $\varphi(x,y) = \varphi 2_1$      $\varphi(x,y) = \varphi 2_2$

Figure 1.1:

$\varphi(x,y) = \varphi2_3$ $\qquad\qquad$ $\varphi(x,y) = \varphi3_1$ $\qquad\qquad$ $\varphi(x,y) = \varphi3_2$

Figure 1.2:



$\varphi(x,y) = \varphi4$ $\qquad\qquad$ $\varphi(x,y) = \varphi5$ $\qquad\qquad$ $\varphi(x,y) = \varphi6$

Figure 1.3:
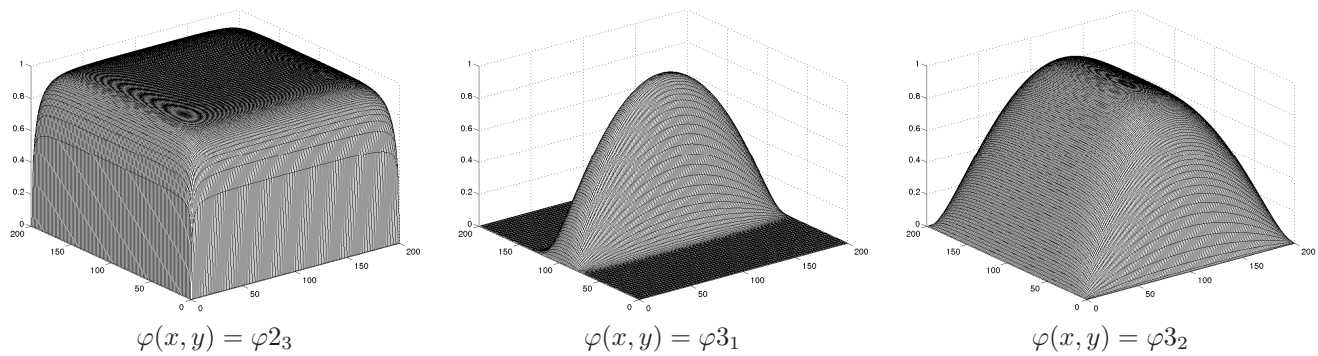


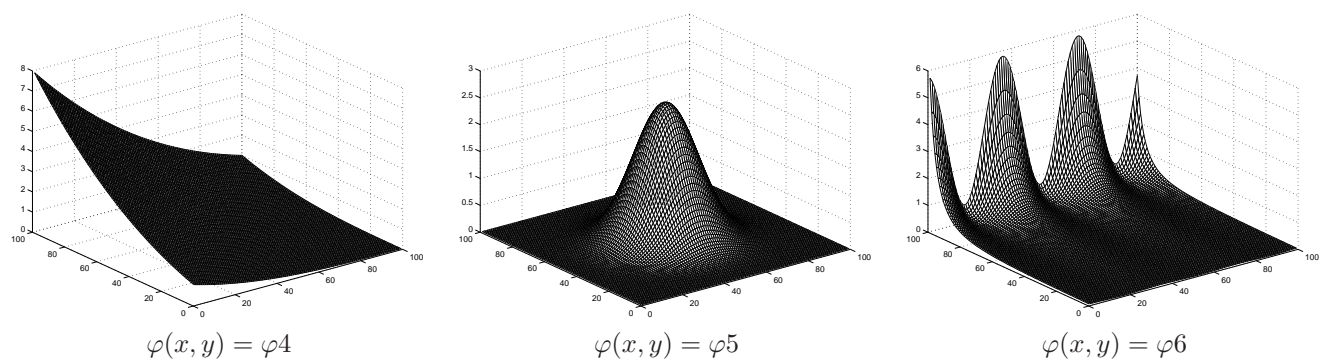$\varphi(x,y) = \varphi7$ $\qquad\qquad$ $\varphi(x,y) = \varphi8$
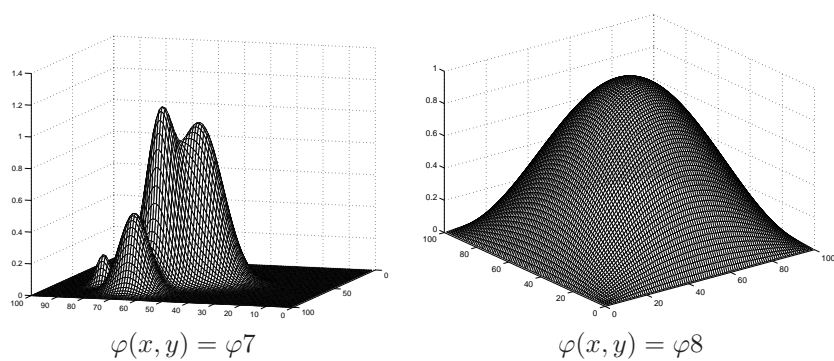
Figure 1.4:

# Chapter 2

## 2.1  The lagged diffusivity functional iteration procedure

Consider a nonlinear steady state reaction diffusion equation of the form

$$-\text{div}(\sigma\nabla\varphi) + \tilde{\boldsymbol{v}}\cdot\nabla\varphi + \alpha\varphi + g(x,y,\varphi) = s(x,y) \tag{2.1}$$

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

In the equation (2.1) the convection term $\tilde{\boldsymbol{v}}\cdot\nabla\varphi$ has been taken into account; however, we will consider only problems for which convection is not dominated.

Equation (2.1) is supplemented by suitable boundary conditions on the contour $\Gamma$ of $\Omega$. We assume that (2.1) has an isolated solution.

When we use a finite difference discretization, the elliptic equation (2.1) supplemented by a Dirichlet boundary condition can be transcribed into a *strongly nonlinear* system of algebraic equations of the form

$$\boldsymbol{F}(\boldsymbol{u}) \equiv A(\boldsymbol{u})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0} \tag{2.2}$$

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 $\boldsymbol{u}$.

Here, the vector $\boldsymbol{u} \in \mathbb{R}^n$ is the approximation of $\varphi(x,y)$ on a mesh of points $\Omega_h$ superimposed on the domain $\Omega$.

We will investigate the solvability of the system of the nonlinear partial difference equation (2.2). We assume that this system has a solution $\boldsymbol{u}^*$.

For solving system (2.2) the easiest and maybe the most common method is *to lag* part of the nonlinear terms in (2.2) (see [39]).

Our purpose here is to re-examine the *diffusivity lagged functional iteration* (LDFI)–procedure for solving the system of nonlinear partial difference equations of elliptic type (2.2) in the context of Parallel Computing.

With this iterative procedure the nonlinear system (2.2) can be solved via a sequence of systems of *weakly nonlinear* difference equations where only $\boldsymbol{G}$, but not $\sigma$ and $\tilde{\boldsymbol{v}}$ in (2.1) depends on the approximate solution $\boldsymbol{u}$ of $\varphi$.

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 (2.2), 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

$$\boldsymbol{F}_\nu(\boldsymbol{u}) \equiv A(\boldsymbol{u}^{(\nu)})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0} \tag{2.3}$$

with the criterion for acceptability of the solution

$$\|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \leq \varepsilon_{\nu+1} \tag{2.4}$$

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 (2.3). This solver must be particularly well suited for implementation on parallel computer.

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

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

$$\tag{2.5}$$

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

where $\tau_1$ and $\tau_2$ are prespecified tolerances.

The purpose of the following three sections is to analyse the convergence of the Lagged Diffusivity Functional Iteration (LDFI) procedure, considering the essential features of the system of nonlinear equations generated by a finite difference discretization of a reaction–diffusion *model problem*.

It is important to define for this model problem the set $\mathcal{B}$ of all grid functions which contains the solutions of the systems and all iterates $\{\boldsymbol{u}^{(\nu)}\}$ of the LDFI–procedure.

For example, the model problem studied in [28] allows to present a helpful paradigm for proving the convergence of the LDFI–procedure.

Here, we will follow this scheme.

At the beginning we summarize some properties of finite difference operators defined in $\mathcal{B}$. Then, we consider the solutions $\{\boldsymbol{u}^{(\nu)}\}$ of the sequence of systems of weakly nonlinear difference equations generated by introducing diffusivity lagging in the original nonlinear system. This solution $\{\boldsymbol{u}^{(\nu)}\}$ can be computed iteratively by the simplified version of the Newton–Arithmetic Mean method for solving weakly nonlinear systems that is particularly suited for implementation on parallel computers ( [9], [10]; see also [2], [3], [4], [41]).

Finally, we prove the convergence of these iterates $\{\boldsymbol{u}^{(\nu)}\}$ to a solution $\boldsymbol{u}^*$ of the original nonlinear system, using well known standard techniques.

In the section of the numerical experiments 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 convection term is present.

## 2.2  A model problem

In a rectangular domain $\Omega$ with boundary $\Gamma$, we consider the reaction diffusion convection equation

$$-\mathrm{div}(\sigma(\boldsymbol{x}, \varphi)\nabla\varphi) + \tilde{\boldsymbol{v}} \cdot \nabla\varphi + \alpha(\boldsymbol{x})\varphi + g(\boldsymbol{x}, \varphi) = s(\boldsymbol{x}) \qquad \boldsymbol{x} \in \Omega \tag{2.6}$$

subject to the homogeneous Dirichlet boundary conditions

$$\varphi(\boldsymbol{x}) = 0 \qquad \boldsymbol{x} \in \Gamma \tag{2.7}$$

Here, $\tilde{\boldsymbol{v}} = (\tilde{v}_1, \tilde{v}_2)^T$ is the velocity vector; $\tilde{v}_1$ and $\tilde{v}_2$ are constants.

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[1] mapping in $\varphi$ with constant $c > 0$ (uniformly in $\boldsymbol{x}$) and is continuously differentiable in $\varphi$.

There exist various techniques for discretizing the problem (2.6)–(2.7). Using the Taylor series approach, equation (2.6) 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 (2.6) 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 [8].

Providing a discretization error $O(h^2)$, the finite difference approximation of (2.6) 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

$$-(B_{ij} + \tilde{B}_{ij})u_{ij-1} - (L_{ij} + \tilde{L}_{ij})u_{i-1j} + (D_{ij} + \tilde{D}_{ij})u_{ij} - \qquad (2.8)$$
$$-(R_{ij} + \tilde{R}_{ij})u_{i+1j} - (T_{ij} + \tilde{T}_{ij})u_{ij+1} + g(x_i, y_j, u_{ij}) - s(x_i, y_j) = 0$$

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

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

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

$$\tilde{L}_{ij} = \frac{\tilde{v}_1}{2h} \qquad \tilde{B}_{ij} = \frac{\tilde{v}_2}{2h} \qquad\qquad (2.9)$$

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

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

---

[1][33, p. 141] A mapping $F : D \subset \mathbb{R}^n \to \mathbb{R}^n$ is uniformly monotone if there exists a constant $\gamma > 0$ such that for all $\boldsymbol{u}, \boldsymbol{v} \in D$ we have

$$(\boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}))^T(\boldsymbol{u} - \boldsymbol{v}) \geq \gamma(\boldsymbol{u} - \boldsymbol{v})^T(\boldsymbol{u} - \boldsymbol{v})$$

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: $k = (j-1) \cdot N + i$. We set $n = N \cdot M$, 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 (2.8) for $i = 1, ..., N$, $j = 1, ..., M$, can be written as formula (2.2)

$$\boldsymbol{F}(\boldsymbol{u}) \equiv A(\boldsymbol{u})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0} \tag{2.10}$$

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 $k = (j-1) \cdot N + i$ (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 $k = (j-1) \cdot N + i$ (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 ([40, p. 23]) and has positive diagonal elements, $a_{ii}(\boldsymbol{u}) > 0$ and nonpositive off diagonal elements $a_{ij}(\boldsymbol{u}) \le 0, i \ne j$, with $i, j = 1, ..., n$; therefore $A(\boldsymbol{u})$ is an M–matrix ([40, p. 91]).

In the case of reaction diffusion equation ($\tilde{\boldsymbol{v}} = 0$), the matrix $A(\boldsymbol{u})$ is also symmetric, then $A(\boldsymbol{u})$ is symmetric positive definite (Stieltjes matrix [40, 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 $k = (j-1) \cdot N + i$. We observe, that $G_k(\boldsymbol{u})$, the $k$–th component of $\boldsymbol{G}(\boldsymbol{u})$, 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 (2.10) is the null vector since we have the homogeneous Dirichlet condition (2.7) 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

$$
\begin{aligned}
< \boldsymbol{u}, \boldsymbol{v} > &= h^2 \sum_{i=1}^{N} \sum_{j=1}^{M} u_{ij} v_{ij} \\[2mm]
\|\boldsymbol{u}\|_h &= \left( h^2 \sum_{i=1}^{N} \sum_{j=1}^{M} |u_{ij}|^2 \right)^{1/2} = (< \boldsymbol{u}, \boldsymbol{u} >)^{1/2}
\end{aligned}
\tag{2.11}
$$

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

$$\|\boldsymbol{u}\|_h \leq \rho \tag{2.12}$$

$$|\nabla_x u_{ij}| \leq \beta \qquad |\nabla_y u_{ij}| \leq \beta \tag{2.13}$$

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 (2.10) 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., [28], where a proof of the existence of such a solution $\boldsymbol{u}^*$ of (2.10) in $\mathcal{B}$ is given; also a condition for which $\boldsymbol{u}^*$ is unique in $\mathcal{B}$ has been obtained)

We wish to compute a solution of the system of nonlinear equations (2.10) with a common iterative procedure in which the nonlinear diffusion part in (2.10) may be evaluated for the previous iteration. This approach is called *nonlinearity lagging* in the diffusivity term.

## 2.3  Some properties of finite difference operators

In the following we summarize some properties of finite difference operators.

**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 (2.9), $L_{ij}$, $R_{ij}$, $B_{ij}$ and $T_{ij}$, are dependent on the grid function $w_{ij}$, then,

$$\sum_{i=1}^{N} \left[L_{ij}(u_{ij} - u_{i-1j}) - R_{ij}(u_{i+1j} - u_{ij})\right] v_{ij} = \tag{2.14}$$

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

and

$$\sum_{j=1}^{M} \left[B_{ij}(u_{ij} - u_{ij-1}) - T_{ij}(u_{ij+1} - u_{ij})\right] v_{ij} = \tag{2.15}$$

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

**Proof.** We prove fomula (2.14). We have[2]

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

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

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

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

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

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

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

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

then, since $v_{0j} = 0$ for (2.7), the expression on the right hand side becomes

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

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

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

and by $u_{N+1j} = 0$, we have formula (2.14). Similarly, we obtain formula (2.15).     ♯

We remark that from the right hand side of (2.14) and (2.15) 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 (2.10), replacing $\boldsymbol{u}$ with $\boldsymbol{w}$, when $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ and $\alpha(x,y) = 0$ in (2.6):

$$< A(\boldsymbol{w})\boldsymbol{u}, \boldsymbol{v} > \quad = \quad h^2 \sum_{j=1}^{M}\sum_{i=1}^{N} \sigma(w_{ij})(\nabla_x u_{ij}\nabla_x v_{ij} + \nabla_y u_{ij}\nabla_y v_{ij}) + \qquad (2.16)$$

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

---

[2]For simplicity, we now omit the coordinates $x$ and $y$ in the expression of the function $\sigma$.

**Proof.** Suppose the coefficients in (2.9), $L_{ij}$, $R_{ij}$, $B_{ij}$ and $T_{ij}$, are functions of the grid function $w_{ij}$, we have,

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

Using formulae (2.14) and (2.15) and keeping into account of the inner product in $l_2(\Omega_h)$, we have

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

Then we have formula (2.16). ♯

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 (2.16) 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 (2.10) and let $A(\boldsymbol{w})$ the matrix $n \times n$ in (2.10) with $\boldsymbol{u}$ replaced by the vector $\boldsymbol{w}$, when $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ and $\alpha(x, y) = 0$ in (2.6).

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

$$
\begin{aligned}
| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \; &\leq \; \frac{h^2 \Lambda \beta \phi}{2} \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right] + \\
&\quad + \frac{\Lambda \beta}{\phi} \| \boldsymbol{u} - \boldsymbol{w} \|_h^2
\end{aligned} \tag{2.17}
$$

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 (2.16) in Lemma 2, we can write

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

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,

$$
\begin{aligned}
| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \quad \leq \quad & \sum_{j=1}^{M}\sum_{i=1}^{N} \left[ h^2 |\sigma(u_{ij}) - \sigma(w_{ij})| \left( |\nabla_x v_{ij}| \cdot |\nabla_x (u_{ij} - v_{ij})| + \right. \right. \\
& \left. \left. + |\nabla_y v_{ij}| \cdot |\nabla_y (u_{ij} - v_{ij})| \right) \right] + \\
& + \sum_{j=1}^{M} |\sigma(u_{N+1j}) - \sigma(w_{N+1j})| \cdot |v_{Nj}| \cdot |u_{Nj} - v_{Nj}| + \\
& + \sum_{i=1}^{N} |\sigma(u_{iM+1}) - \sigma(w_{iM+1})| \cdot |v_{iM}| \cdot |u_{iM} - v_{iM}|
\end{aligned}
\qquad (2.18)
$$

The assumption (iii) implies that, for a 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 $j = 1, ..., M + 1$

$$
|\sigma(u_{ij}) - \sigma(w_{ij})| \leq \Lambda |u_{ij} - w_{ij}| \qquad (2.19)
$$

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

Furthermore, Property A assures that there exists a constant $\beta > 0$ such that inequality (2.13) holds

$$
|\nabla_x v_{ij}| \leq \beta \qquad |\nabla_y v_{ij}| \leq \beta
$$

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 inqualities (2.19) and (2.13) into the espression in (2.18) we obtain

$$
\begin{aligned}
| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \quad \leq \quad & \sum_{j=1}^{M}\sum_{i=1}^{N} \left[ h^2 \Lambda \beta |u_{ij} - w_{ij}| \left( |\nabla_x (u_{ij} - v_{ij})| + |\nabla_y (u_{ij} - v_{ij})| \right) \right] + \\
& + \Lambda \sum_{j=1}^{M} |u_{N+1j} - w_{N+1j}| \cdot |v_{Nj}| \cdot |u_{Nj} - v_{Nj}| + \\
& + \Lambda \sum_{i=1}^{N} |u_{iM+1} - w_{iM+1}| \cdot |v_{iM}| \cdot |u_{iM} - v_{iM}|
\end{aligned}
$$

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

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

The last expression can be written

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

Using a well known technical trick, i.e.

$$\sqrt{ab} \leq \frac{1}{2}a + \frac{1}{2}b \qquad a > 0, b > 0$$

we have

$$| < (A(\boldsymbol{u}) - A(\boldsymbol{w}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > | \quad \leq \quad h^2 \Lambda \beta \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ \frac{|u_{ij} - w_{ij}|^2}{2\phi} + |\nabla_x(u_{ij} - v_{ij})|^2 \frac{\phi}{2} + \right.$$
$$\left. + \frac{|u_{ij} - w_{ij}|^2}{2\phi} + |\nabla_y(u_{ij} - v_{ij})|^2 \frac{\phi}{2} \right] \tag{2.20}$$

and we obtain formula (2.17).          ♯

As consequence of lemmas 1, 2 and 3 we prove the corollary on the uniform monotonicity of the mapping $\boldsymbol{F}(\boldsymbol{u})$. Thus, from Hadamard Theorem ([20]), the nonlinear system (2.10) has a unique solution (e.g. [33, p. 143]).

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 (2.2) exists and is unique; besides, it is possible to construct an iterative procedure that can guarantee a global convergence to the solution of (2.2) ([27]).

**Corollary 1.** From lemmas 1, 2, and 3 we have that the mapping $\boldsymbol{F}(\boldsymbol{u})$ is uniformly monotone in $\mathcal{B}$.

**Proof.** Indeed, we prove that, under suitable conditions (i)–(iv), the mapping $\boldsymbol{F}(\boldsymbol{u})$ in (2.10), with $\tilde{\boldsymbol{v}} = \boldsymbol{0}$ and $\alpha(x, y) = 0$ in (2.6), is uniformly monotone on $\mathcal{B}$.

By definition of uniformly monotone mapping, we should prove it exists a positive constant $\gamma$ such that

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad \gamma < \boldsymbol{u} - \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > \tag{2.21}$$

From

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

we have

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad = \quad < \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{G}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > + \tag{2.22}$$
$$+ < A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > +$$
$$+ < (A(\boldsymbol{u}) - A(\boldsymbol{v}))\,\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >$$

We separately examine the three terms of the right hand side of (2.22).

Since $g$ is uniformly monotone respect to $\varphi$, and $\boldsymbol{G}(\boldsymbol{u})$ is a diagonal mapping, there exists a positive constant $c$ such that, for all $\boldsymbol{u}, \boldsymbol{v} \in \mathcal{B}$,

$$(G_i(u_i) - G_i(v_i)) \cdot (u_i - v_i) \geq c(u_i - v_i)^2$$

for $i = 1, ..., n$. Thus for the discrete $l^2(\Omega_h)$ inner product (2.11) we have

$$< \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{G}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad c < \boldsymbol{u} - \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >= c\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 \tag{2.23}$$

By using Lemma 2 we can write

$$< A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad = \quad h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \sigma(u_{ij})[(\nabla_x(u_{ij} - v_{ij}))^2 + (\nabla_y(u_{ij} - v_{ij}))^2] +$$
$$+ \sum_{j=1}^{M} \sigma(u_{N+1j})(u_{Nj} - v_{Nj})^2 + \sum_{i=1}^{N} \sigma(u_{iM+1})(u_{iM} - v_{iM})^2$$

The assumption (ii) on the the uniform lower boundedness of $\sigma$ respect to the variable $\varphi$ and the positivity of terms to the square, permits to write

$$< A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad h^2 \sigma_{\min} \sum_{j=1}^{M} \sum_{i=1}^{N} \left( |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right) +$$
$$+ \sigma_{\min} \left[ \sum_{j=1}^{M} |u_{Nj} - v_{Nj}|^2 + \sum_{i=1}^{N} |u_{iM} - v_{iM}|^2 \right] \tag{2.24}$$

By using Lemma 3 we have

$$< (A(\boldsymbol{u}) - A(\boldsymbol{v}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad -\frac{h^2 \Lambda \beta \phi}{2} \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right] -$$
$$- \frac{\Lambda \beta}{\phi} \|\boldsymbol{u} - \boldsymbol{v}\|_h^2 \tag{2.25}$$

By formulae (2.23), (2.25) and (2.24), the formula (2.22) can be written

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad c\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 - \frac{\Lambda \beta}{\phi} \|\boldsymbol{u} - \boldsymbol{v}\|_h^2 - \tag{2.26}$$
$$- \frac{h^2 \Lambda \beta \phi}{2} \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right] +$$
$$+ h^2 \sigma_{\min} \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right] +$$
$$+ \sigma_{\min} \left[ \sum_{j=1}^{M} |u_{Nj} - v_{Nj}|^2 + \sum_{i=1}^{N} |u_{iM} - v_{iM}|^2 \right]$$

The last two terms of the previous expression are positive, then we can write

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad (c - \frac{\Lambda \beta}{\phi})\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 +$$
$$+ (\sigma_{\min} - \frac{\Lambda \beta \phi}{2})h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \left[ |\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2 \right]$$

If we set

$$\phi = \frac{2\sigma_{\min}}{\Lambda \beta}$$

we have

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad (c - \frac{\Lambda^2 \beta^2}{2\sigma_{\min}})\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 \tag{2.27}$$

When

$$c - \frac{\Lambda^2 \beta^2}{2\sigma_{\min}} > 0 \tag{2.28}$$

or

$$\frac{\Lambda^2 \beta^2}{2c\sigma_{\min}} < 1 \qquad (2.29)$$

inequality (2.27) shows that the mapping $\boldsymbol{F}(\boldsymbol{u})$ is uniformly monotone on $\mathcal{B}$ with constant $\gamma \equiv c - (\Lambda^2\beta^2)/(2\sigma_{\min})$ in (2.21).

When inequality (2.27) holds, there exists only one solution $\boldsymbol{u}^*$ of the system $\boldsymbol{F}(\boldsymbol{u}) = 0$ in $\mathcal{B}$. Indeed, suppose that $\tilde{\boldsymbol{u}}$ is some other solution of $\boldsymbol{F}(\boldsymbol{u}) = \boldsymbol{0}$ in $\mathcal{B}$ and $\tilde{\boldsymbol{u}} \neq \boldsymbol{u}^*$. Then, since the mapping $\boldsymbol{F}(\boldsymbol{u})$ is uniformly monotone on $\mathcal{B}$ and (2.27) and (2.28) hold, we have

$$0 = < \boldsymbol{F}(\boldsymbol{u}^*) - \boldsymbol{F}(\tilde{\boldsymbol{u}}), \boldsymbol{u}^* - \tilde{\boldsymbol{u}} > \geq (c - \frac{\Lambda^2\beta^2}{2\sigma_{\min}})\|\boldsymbol{u}^* - \tilde{\boldsymbol{u}}\|_h^2 > 0$$

which is a contradiction. Hence $\boldsymbol{u}^*$ is the only solution of $\boldsymbol{F}(\boldsymbol{u}) = 0$ in $\mathcal{B}$.

For sake of completeness, we show the uniform monotonicity of $\boldsymbol{F}(\boldsymbol{u})$ also in the cases $\alpha(x, y) > 0$ and $\tilde{\boldsymbol{v}} \neq \boldsymbol{0}$ for the problem (2.6).

When we suppose $\alpha(x, y) > 0$ in (2.6), we can write $A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{D}$ then it is easy to show that formula (2.16) in Lemma 2 becomes

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

where $\tilde{D}_{ij}$ are the diagonal elements of the matrix $\tilde{D}$ as in (2.9), and since

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

formula (2.17) in Lemma 3 is unchanged. Furthermore the formula (2.22) becomes

$$
\begin{aligned}
< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad = \quad & < \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{G}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > + \\
& + < A_1(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > + \\
& + < (A_1(\boldsymbol{u}) - A_1(\boldsymbol{v}))\boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > + \\
& + < \tilde{D}(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} >
\end{aligned}
$$

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

$$< \tilde{D}(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > = h^2\sum_{j=1}^{M}\sum_{i=1}^{N}\alpha(x_i, y_j)(u_{ij} - v_{ij})(u_{ij} - v_{ij}) \geq \alpha_{\min}\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 \qquad (2.30)$$

thus, formula (2.26) is rewritten as

$$
\begin{aligned}
< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad & c\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 + \alpha_{\min}\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 - \frac{\Lambda\beta}{\phi}\|\boldsymbol{u} - \boldsymbol{v}\|_h^2 - \\
& - \frac{h^2\Lambda\beta\phi}{2}\sum_{j=1}^{M}\sum_{i=1}^{N}\left[|\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2\right] + \\
& + h^2\sigma_{\min}\sum_{j=1}^{M}\sum_{i=1}^{N}\left[|\nabla_x(u_{ij} - v_{ij})|^2 + |\nabla_y(u_{ij} - v_{ij})|^2\right] + \\
& + \sigma_{\min}\sum_{j=1}^{M}\sum_{i=1}^{N}\left[|u_{Nj} - v_{Nj}|^2 + |u_{iM} - v_{iM}|^2\right]
\end{aligned}
$$

Therefore, when

$$c + \alpha_{\min} - \frac{\Lambda^2 \beta^2}{2\sigma_{\min}} > 0$$

inequality (2.27) becomes

$$< \boldsymbol{F}(\boldsymbol{u}) - \boldsymbol{F}(\boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \quad \geq \quad (c + \alpha_{\min} - \frac{\Lambda^2 \beta^2}{2\sigma_{\min}}) \|\boldsymbol{u} - \boldsymbol{v}\|_h^2$$

and guarantees the uniform monotonicity of $\boldsymbol{F}(\boldsymbol{u})$ in $\mathcal{B}$.

When we suppose $\tilde{\boldsymbol{v}} \neq \boldsymbol{0}$ in (2.6), we can write $A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{A}$ with $\tilde{A}$ a skew–symmetric matrix (i.e. $\tilde{A} = -\tilde{A}^T$).
We denote with $A_s$ the symmetric part of $A(\boldsymbol{u})$ and $A_{ss}$ the skew–symmetric part of $A(\boldsymbol{u})$, thus

$$A_s = \frac{A(\boldsymbol{u}) + A(\boldsymbol{u})^T}{2} = A_1(\boldsymbol{u}) \qquad A_{ss} = \frac{A(\boldsymbol{u}) - A(\boldsymbol{u})^T}{2} = \tilde{A}$$

Since, for all vector $\boldsymbol{w} \in \mathbb{R}^n$ we have

$$\boldsymbol{w}^T A(\boldsymbol{u})\boldsymbol{w} = \boldsymbol{w}^T A_s \boldsymbol{w}$$

and for all $\boldsymbol{u}, \boldsymbol{v} \in \mathcal{B}$, we have

$$A(\boldsymbol{u}) - A(\boldsymbol{v}) = (A_1(\boldsymbol{u}) + \tilde{A}) - (A_1(\boldsymbol{v}) + \tilde{A}) = A_1(\boldsymbol{u}) - A_1(\boldsymbol{v})$$

then, we otain

$$< A(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} >=< A_1(\boldsymbol{u})(\boldsymbol{u} - \boldsymbol{v}), \boldsymbol{u} - \boldsymbol{v} > \tag{2.31}$$

and

$$< (A(\boldsymbol{u}) - A(\boldsymbol{v}))\, \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} >=< (A_1(\boldsymbol{u}) - A_1(\boldsymbol{v}))\, \boldsymbol{v}, \boldsymbol{u} - \boldsymbol{v} > \tag{2.32}$$

Thus, Lemma 3 holds also for $A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{A}$ and formula (2.22) and inequality (2.27) with condition (2.28) guarantee the uniform monotonicity of $\boldsymbol{F}(\boldsymbol{u})$ in $\mathcal{B}$.      ♯

## 2.4   Convergence of the LDFI–procedure

We will now investigate the solvability of the system of nonlinear difference equations (2.10) by applying the LDFI–procedure.
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 (2.6)–(2.7) can be solved via a sequence of systems of weakly nonlinear difference equations where only $\boldsymbol{G}$ but not $\sigma$ depend on the approximate solution $\boldsymbol{u}$ of $\varphi$.
Specifically, if $\boldsymbol{u}^{(\nu)}$ is an estimate of the solution $\boldsymbol{u}^*$ of (2.10), we will determine a new estimate of $\boldsymbol{u}^*$ by solving the weakly nonlinear system (2.3)

$$\boldsymbol{F}_\nu(\boldsymbol{u}) \equiv A(\boldsymbol{u}^{(\nu)})\boldsymbol{u} + \boldsymbol{G}(\boldsymbol{u}) - \boldsymbol{s} = \boldsymbol{0} \tag{2.33}$$

An approximate solution of the weakly nonlinear system (2.33) is computed by the simplified Newton Arithmetic Mean 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

$$\|\boldsymbol{F}_\nu(\boldsymbol{u}^{(\nu+1)})\| \leq \varepsilon_{\nu+1} \tag{2.34}$$

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 preavious 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 (2.13), 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

$$|\nabla_x u_{ij}^{(\nu)}| \leq \beta + \varepsilon_\nu \beta_0 \qquad |\nabla_y u_{ij}^{(\nu)}| \leq \beta + \varepsilon_\nu \beta_0 \tag{2.35}$$

instead of (2.13), $i = 1, ..., N + 1$ 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 (2.6)–(2.7) subject to the conditions (i)–(iv) and $A(\boldsymbol{u})$ the irreducibly M–matrix in (2.10).
Let $\{\boldsymbol{u}^{(\nu+1)}\}$, $\nu = 0, 1, ...$, be the solution of $\boldsymbol{F}_\nu(\boldsymbol{u}) = \boldsymbol{0}$ satisfying the condition (2.34) with $\boldsymbol{F}_\nu(\boldsymbol{u})$ as in (2.33) and satisfying Property A with (2.35) instead of (2.13).
Then, the sequence $\{\boldsymbol{u}^{(\nu)}\}$ converges to $\boldsymbol{u}^*$.

**Proof.** First we consider the case of $\alpha(x, y) > 0$ and $\tilde{\boldsymbol{v}} \neq \boldsymbol{0}$ for the problem (2.6)–(2.7). The solution $\boldsymbol{u}^*$ in $\mathcal{B}$ of (2.10) 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 (2.34).
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

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

$$\tag{2.36}$$

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

Using (2.16) and assumption (ii), we can write

$$
\begin{aligned}
< A(\boldsymbol{u}^*)(\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > \quad = \quad & h^2 \sum_{j=1}^{M} \sum_{i=1}^{N} \sigma(u_{ij}^*) \cdot \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) + \\
& + \sum_{j=1}^{M} \sigma(u_{N+1j}^*) \cdot |u_{Nj}^* - u_{Nj}^{(\nu+1)}|^2 + \sum_{i=1}^{N} \sigma(u_{iM+1}^*) \cdot |u_{iM}^* - u_{iM}^{(\nu+1)}|^2 \\
\geq \quad & \sigma_{\min} \sum_{j=1}^{M} \sum_{i=1}^{N} h^2 \left( |\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 + |\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) \\
& + \sigma_{\min} \left( \sum_{j=1}^{M} |u_{Nj}^* - u_{Nj}^{(\nu+1)}|^2 + \sum_{i=1}^{N} |u_{iM}^* - u_{iM}^{(\nu+1)}|^2 \right)
\end{aligned}
$$

$$\tag{2.37}$$

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

$$(G_{ij}(u_{ij}) - G_{ij}(v_{ij})) \cdot (u_{ij} - v_{ij}) \geq c(u_{ij} - v_{ij})^2$$

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 (2.11) we have the inequality

$$< \boldsymbol{G}(\boldsymbol{u}^*) - \boldsymbol{G}(\boldsymbol{u}^{(\nu+1)}), \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > \quad \geq \quad c\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2 \tag{2.38}$$

Using Lemma 3 (see formula (2.20) 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 (2.35), we can write

$$
\begin{aligned}
| < (A(\boldsymbol{u}^*) - A(\boldsymbol{u}^{(\nu)}))\boldsymbol{u}^{(\nu+1)}, \boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)} > | \quad \leq \quad & \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)}{2} \cdot \\
& \cdot \sum_{j=1}^{M}\sum_{i=1}^{N} h^2 \left( \frac{|u_{ij}^* - u_{ij}^{(\nu)}|^2}{\phi} + \phi|\nabla_x(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right) + \\
& + \frac{\Lambda(\beta + \varepsilon_{\nu+1}\beta_0)}{2} \cdot \\
& \cdot \sum_{j=1}^{M}\sum_{i=1}^{N} h^2 \left( \frac{|u_{ij}^* - u_{ij}^{(\nu)}|^2}{\phi} + \phi|\nabla_y(u_{ij}^* - u_{ij}^{(\nu+1)})|^2 \right)
\end{aligned}
\tag{2.39}
$$

It now follows from (2.36) that

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

and from (2.37), (2.38) and (2.39) that

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

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

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

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

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

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

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

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

we obtain

$$
c\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2 - \frac{\Lambda^2(\beta + \varepsilon_{\nu+1}\beta_0)^2}{2\sigma_{\min}}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2 \leq \varepsilon_{\nu+1}\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h
\tag{2.40}
$$

Since then 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 (2.40) we have the inequality

$$
\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu+1)}\|_h^2 \leq \gamma\|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2 + a\varepsilon_{\nu+1}
\tag{2.41}
$$

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 [28], if there exists an integer $\nu_0$ such that $\gamma < 1$ for all $\nu \geq \nu_0$, we can write (2.41) 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 ([33, p. 399], [42, 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 (2.10).

For sake of completeness, it easy to show that we have the convergence of $\{\boldsymbol{u}^{(\nu)}\}$ to the solution $\boldsymbol{u}^*$ of the system (2.10) also in the cases $\alpha(x, y) > 0$ and $\tilde{\boldsymbol{v}} \neq \boldsymbol{0}$ for the problem (2.6)–(2.7).
Indeed, since for $A(\boldsymbol{u}) = A_1(\boldsymbol{u}) + \tilde{A} + \tilde{D}$ formulae (2.31) and (2.32) hold, formula (2.36) becomes

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

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

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

Since (2.30) then, at the left hand side of inequality (2.40) we have to add the term $\alpha_{\min} \|\boldsymbol{u}^* - \boldsymbol{u}^{(\nu)}\|_h^2$ and in (2.41) the parameter $\gamma$ becomes

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

## 2.5   Solution of the weakly nonlinear system

In order to define the inner iterative solver for the nonlinear system (2.3) (or (2.33)), setting $\boldsymbol{w}^{(0)} = \boldsymbol{u}^{(\nu)}$, the simplified–Newton method finds the solution $\Delta \boldsymbol{w}^{(k)}$ of

$$C_\nu \Delta \boldsymbol{w} = -\boldsymbol{F}_\nu(\boldsymbol{w}^{(k)}) \tag{2.42}$$

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

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \Delta \boldsymbol{w}^{(k)} \tag{2.43}$$

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 (2.42)–(2.43) are rewritten in such a way that the vector $\boldsymbol{w}^{(k+1)}$ is the solution of the linear system

$$C_\nu \boldsymbol{w} = G'(\boldsymbol{u}^{(\nu)})\boldsymbol{w}^{(k)} - \boldsymbol{G}(\boldsymbol{w}^{(k)}) + \boldsymbol{s} \tag{2.44}$$

for $k = 0, 1, ....$
The system (2.44) is solved by the block version of the Arithmetic Mean (AM) method introduced in ([14]).

We consider the following decomposition of the matrix

$$C_\nu = \begin{pmatrix} C_{11}(\boldsymbol{u}^{(\nu)}) & C_{12}(\boldsymbol{u}^{(\nu)}) & & & \\ C_{21}(\boldsymbol{u}^{(\nu)}) & C_{22}(\boldsymbol{u}^{(\nu)}) & C_{23}(\boldsymbol{u}^{(\nu)}) & & \\ & & \ddots & & \ddots \\ & & & C_{MM-1}(\boldsymbol{u}^{(\nu)}) & C_{MM}(\boldsymbol{u}^{(\nu)}) \end{pmatrix} \tag{2.45}$$

into the two splittings

$$C_\nu = H_1(\boldsymbol{u}^{(\nu)}) - K_1(\boldsymbol{u}^{(\nu)}) = H_2(\boldsymbol{u}^{(\nu)}) - K_2(\boldsymbol{u}^{(\nu)}) \tag{2.46}$$

where, if $M$ is even

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

and, consequently

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

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

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-Arithmetic Mean 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

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

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

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

$\boldsymbol{w}^{(k+1)} = \boldsymbol{z}_k^{(j_k)}$
$\tag{2.47}$

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

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

(see formula (2.4)).

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 [14], [15], [16], [18], [19].

From hypotheses (i)–(iv), for any vector $\boldsymbol{u}$ and $\boldsymbol{u}^{(\nu)}$, the following **Standard Assumptions** are satisfied:

- $A(\boldsymbol{u}^{(\nu)})$ 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 non-positive off-diagonal entries for all the mesh spacings sufficiently small and for all the iterates $\boldsymbol{u}^{(\nu)} \in \mathbb{R}^n$.

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

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 $\boldsymbol{u} \in \mathbb{R}^n$ and for all the iterates $\boldsymbol{u}^{(\nu)}$ (see, e.g., [32, p. 109]).

We report a general result on the convergence of the simplified Newton–Arithmetic Mean method when the Standard Assumptions are satisfied.

First we should define the matrix ($\rho \geq 0$)

$$M_\nu^{-1} = \frac{1}{2}[(H_1(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1} + (H_2(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}] \tag{2.48}$$

and the iteration matrix

$$H_\nu = \frac{1}{2}[(H_1(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}(K_1(\boldsymbol{u}^{(\nu)}) + \rho I) + (H_2(\boldsymbol{u}^{(\nu)}) + \rho I)^{-1}(K_2(\boldsymbol{u}^{(\nu)}) + \rho I)] \tag{2.49}$$

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

**Theorem 2.** Suppose the system (2.3) $\boldsymbol{F}_\nu(\boldsymbol{u}) = \boldsymbol{0}$ has a solution $\tilde{\boldsymbol{u}} \in \mathbb{R}^n$; assume that Standard Assumptions hold for $\boldsymbol{u} \in \mathbb{R}^n$ (or in an open neighbourhood $\mathcal{K}$ of $\tilde{\boldsymbol{u}}$) and that (2.46), 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 $C_\nu = F'_\nu(\boldsymbol{u}^{(\nu)})$, $\boldsymbol{u}^{(\nu)} \in \mathbb{R}^n$ (or $\boldsymbol{u}^{(\nu)} \in \mathcal{K}$), with the matrix $H_\nu$ in (2.49) convergent.

Then, for any $j_k \geq 1$, the solution $\tilde{\boldsymbol{u}}$ is an attraction point of the simplified Newton–Arithmetic Mean iteration $\{\boldsymbol{w}^{(k)}\}$ defined in (2.47).

**Proof.** The Standard Assumptions assure that the Jacobian matrix $F'_\nu(\boldsymbol{u})$ is continuous and nonsingular and a monotone matrix in $\mathbb{R}^n$ (or 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 (2.48) and (2.49) are nonnegative and $H_\nu$ is a convergent matrix, $\rho(H_\nu) < 1$ ([31]).

Then, from the identity

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

it is possible to write the simplified Newton–Arithmetic Mean iteration (2.47) as

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

that is a *generalized linear iteration* and the proof runs as the one of 10.3.1 in [33, 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 [9], [10], [11], [12], [13].

## 2.6   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 (2.6) with homogeneous Dirichlet boundary conditions (2.7) and with nonhomogeneous Dirichlet boundary conditions

$$\varphi(\boldsymbol{x}) = U_0(\boldsymbol{x}) \qquad \boldsymbol{x} \in \Gamma \tag{2.50}$$

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 (2.2), $i = 1, ..., N$, $j = 1, ..., M$ or a prespecified exact solution $\varphi(x, y)$ of the differential problem (2.6); 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:

$$
\begin{aligned}
\sigma 1 \quad &: \quad \sigma(\varphi) = 1 \\
\sigma 2 \quad &: \quad \sigma(\varphi) = \frac{\varphi^2}{\varphi + 1 + 10^{-2}} + 10^{-2} \\
\sigma 3 \quad &: \quad \sigma(\varphi) = \varphi^2 + 10^{-1} \\
\sigma 4 \quad &: \quad \sigma(\varphi) = \frac{1}{\cosh^2(\varphi)} \\
\sigma 5 \quad &: \quad \sigma(\varphi) = \varphi^{1/2} + 10^{-2} \\
\sigma 6 \quad &: \quad \sigma(\varphi) = \varphi^{1/3} + 10^{-2} \\
\sigma 7_1 \quad &: \quad \sigma(\varphi) = 0.02 + 0.5\phi \\
\sigma 7_2 \quad &: \quad \sigma(\varphi) = 0.5(1 + \phi) \\
\sigma 8 \quad &: \quad \sigma(\varphi) = 0.02 + 0.5\varphi^2 \\
\sigma 9 \quad &: \quad \sigma(\varphi) = 1/(0.02 + 0.5\varphi)
\end{aligned}
$$

Following the notation of Part I, the functions $g$ are dependent on $\varphi$ and are:

$$
\begin{aligned}
g 1_1 \quad &: \quad g(\varphi) = e^{\varphi} \\
g 1_2 \quad &: \quad g(\varphi) = 100 e^{0.5\varphi} \\
g 1_3 \quad &: \quad g(\varphi) = -0.5 e^{\varphi} \\
g 2_1 \quad &: \quad g(\varphi) = \frac{\varphi}{(1 + \varphi)} \\
g 2_2 \quad &: \quad g(\varphi) = \frac{1000\varphi}{(1 + 10\varphi)} \\
g 3 \quad &: \quad g(\varphi) = -(0.4 - \varphi)e^{(-15/(1+\varphi))} \\
g 4 \quad &: \quad g(\varphi) = \frac{0.02\varphi^2}{(3 + \varphi)} \\
g 5_1 \quad &: \quad g(\varphi) = 0.005\varphi \log(1 + \varphi) \\
g 5_2 \quad &: \quad g(\varphi) = 5\varphi \log(1 + \varphi) \\
g 5_3 \quad &: \quad g(\varphi) = 500\varphi \log(1 + \varphi) \\
g 5_4 \quad &: \quad g(\varphi) = 80\varphi \log(1 + \varphi) \\
g 6 \quad &: \quad g(\varphi) = -\varphi(2 - \varphi)
\end{aligned}
$$

Furthermore, we indicate as $g0$ the null function $g = 0$.
We observe that

for $g1_1$, $g1_2$: $g > 0$, $g' > 0$ and $g'' > 0$ for any value of $\varphi$;

for $g2_k$ $(k = 1, 2)$: $g \geq 0$ and $g' > 0$ when $\varphi \geq 0$;

for $g3$: $g \geq 0$, $g' \geq 0$ and $g'' \geq 0$ when $\varphi \geq 0.4$;

for $g4$ and $g5_k$ $(k = 1, 4)$: $g \geq 0$, $g' \geq 0$ and $g'' > 0$ when $\varphi \geq 0$;

and then, for the functions $g1_1$, $g1_2$, $g2_1$, $g2_2$, $g4$ and $g5_k$, $(k = 1, ..., 4)$, the Standard Assumption on $g$ is satisfied for $\varphi \geq 0$.

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

$$\alpha1 \quad : \quad \alpha(x, y) = c(x^3 + y) \qquad \text{with } c = 10, 100, 1000$$
$$\alpha2 \quad : \quad \alpha(x, y) = c\frac{\gamma}{(\varepsilon + x + y)^2} \qquad \gamma = 10 \quad \varepsilon = 10^{-3} \quad \text{and with } c = 1, 10, 100, 1000$$

Now we list the different functions for the exact solution.

$$\varphi1 \quad : \quad \varphi(x, y) = (x - a)(y - a)(b - x)(b - y)(2x^2 + y) \qquad a = 0 \quad b = 1$$
$$\Omega \cup \Gamma = [0, 1] \times [0, 1]$$

Set

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

and

$$\varphi(x, y) = \begin{cases} p(x) \cdot p(y) & 0 < x \leq 1 & 0 < y \leq 1 \\ q(x) \cdot p(y) & 1 < x < 2 & 0 < y \leq 1 \\ p(x) \cdot q(y) & 0 < x \leq 1 & 1 < y < 2 \\ q(x) \cdot q(y) & 1 < x < 2 & 1 < y < 2 \\ 0 & 0 \leq x \leq 2 & y = 0, y = 2 \\ 0 & x = 0, x = 2 & 0 \leq y \leq 2 \end{cases} \tag{2.51}$$

then

$$\varphi2_1 \quad : \quad \varphi(x, y) \text{ as in (2.51)}; \qquad \hat{\alpha} = 100$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2]$$
$$\varphi2_2 \quad : \quad \varphi(x, y) \text{ as in (2.51)}; \qquad \hat{\alpha} = 0.05$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2]$$
$$\varphi2_3 \quad : \quad \varphi(x, y) \text{ as in (2.51)}; \qquad \hat{\alpha} = 0.005$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2]$$

Set

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

and

$$\varphi(x, y) = \begin{cases} p(x) \cdot r(y) & 0 < x \leq 1 & 0 < y < 2 \\ q(x) \cdot r(y) & 1 < x < 2 & 0 < y < 2 \\ 0 & 0 \leq x \leq 2 & y = 0, y = 2 \\ 0 & x = 0, x = 2 & 0 \leq y \leq 2 \end{cases} \tag{2.52}$$

then

$$\varphi3_1 \quad : \quad \varphi(x, y) \text{ as in (2.52)}; \qquad \hat{\alpha} = 100$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2]$$
$$\varphi3_2 \quad : \quad \varphi(x, y) \text{ as in (2.52)}; \qquad \hat{\alpha} = 0.05$$
$$\Omega \cup \Gamma = [0, 2] \times [0, 2]$$

Furthermore we have

$$\begin{aligned}
\varphi 4 \quad &: \quad \varphi(x,y) = (1 + x - y)^3 \\
&\quad \Omega \cup \Gamma = [0,1] \times [0,1] \\
\varphi 5 \quad &: \quad \varphi(x,y) = (8/\pi)e^{-8(x^2+y^2)} \\
&\quad \Omega \cup \Gamma = [-1,1] \times [-1,1] \\
\varphi 6 \quad &: \quad \varphi(x,y) = \frac{3}{2}\left(\cos(\frac{3}{5}(y-1)) + \frac{5}{4}\right) / \left((1 + \frac{x-4}{3})^2\right) e^{-8(x^2+y^2)} \\
&\quad \Omega \cup \Gamma = [-30,-1] \times [-30,-1] \\
\varphi 7 \quad &: \quad \varphi(x,y) = \left(e^{\frac{-(x-3)^2-(y-3)^2}{4}} + e^{-\frac{x^2}{7}-\frac{y^2}{10}} - \frac{1}{5}e^{-(x-5)^2-(y-8)^2} + \frac{1}{2}e^{\frac{-(x-8)^2-(y-4)^2}{4}}\right) \\
&\quad \Omega \cup \Gamma = [-13,13] \times [-13,13] \\
\varphi 8 \quad &: \quad \varphi(x,y) = \sin(\pi x)\sin(\pi y) \\
&\quad \Omega \cup \Gamma = [0,1] \times [0,1]
\end{aligned}$$

We notice that the functions $\varphi 1$, $\varphi 2_1$, $\varphi 2_2$, $\varphi 2_3$, $\varphi 3_1$, $\varphi 3_2$ and $\varphi 8$ are equal to zero at the points of the boundary $\Gamma$.

In the case of nonhomogeneous boundary conditions (2.50) the system (2.3) (or (2.33)) becomes

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

where the vector $\boldsymbol{b}$ depends on the current vector $\boldsymbol{u}^{(\nu)}$ and is

$$\boldsymbol{b} = \begin{pmatrix}
-(B_{11} + \tilde{B}_{11})U_0(x_1,y_0) - (L_{11} + \tilde{L}_{11})U_0(x_0,y_1) \\
-(B_{21} + \tilde{B}_{21})U_0(x_2,y_0) \\
\vdots \\
-(B_{N-11} + \tilde{B}_{N-11})U_0(x_{N-1},y_0) \\
-(B_{N1} + \tilde{B}_{N1})U_0(x_N,y_0) - (R_{N1} + \tilde{R}_{N1})U_0(x_{N+1},y_1) \\
-(L_{12} + \tilde{L}_{12})U_0(x_0,y_2) \\
0 \\
\vdots \\
0 \\
-(R_{N2} + \tilde{R}_{N2})U_0(x_{N+1},y_2) \\
\vdots \\
\vdots \\
-(L_{1M-1} + \tilde{L}_{1M-1})U_0(x_0,y_{M-1}) \\
0 \\
\vdots \\
0 \\
-(R_{NM-1} + \tilde{R}_{NM-1})U_0(x_{N+1},y_{M-1}) \\
-(T_{1M} + \tilde{T}_{1M})U_0(x_1,y_{M+1}) - (L_{1M} + \tilde{L}_{1M})U_0(x_0,y_M) \\
-(T_{2M} + \tilde{T}_{2M})U_0(x_2,y_{M+1}) \\
\vdots \\
-(T_{N-1M} + \tilde{T}_{N-1M})U_0(x_{N-1},y_{M+1}) \\
-(T_{NM} + \tilde{T}_{NM})U_0(x_N,y_{M+1}) - (R_{NM} + \tilde{R}_{NM})U_0(x_{N+1},y_M)
\end{pmatrix}$$

where the terms $B_{ij}$, $L_{ij}$, $R_{ij}$, $T_{ij}$ and $\tilde{B}_{ij}$, $\tilde{L}_{ij}$, $\tilde{R}_{ij}$, $\tilde{T}_{ij}$ are as in (2.9) with $u_{0j}$, $u_{N+1j}$, $u_{i0}$ and $u_{iM+1}$ equal to $U_0(x_0,y_j)$, $U_0(x_{N+1},y_j)$, $U_0(x_i,y_0)$ and $U_0(x_i,y_{M+1})$ respectively, $i = 1,...,N$, $j = 1,...,M$.

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 (2.5)

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

and

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

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 $\varepsilon_1 = 0.1 \cdot \|\boldsymbol{F}(\boldsymbol{u}^{(0)})\|$ and $\varepsilon_{\nu+1} = \min\{0.5 \cdot \varepsilon_\nu, \underline{\varepsilon}\}$, $\nu = 1, 2, ....$ The threshold $\underline{\varepsilon}$ is chosen $10^{-5}$, $10^{-3}$ or $10^{-2}$.

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 $7.02(-10)$ indicates $7.02 \cdot 10^{-10}$.

We indicate with $j_k$ the number of iterations of the Arithmetic Mean method for the solution of the system (2.44)

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

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. Furthermore, we denote as *jtot* the sum of the iterations of the Arithmetic Mean method in the LDFI–procedure.

When we include the discretization error, i.e., the source term is computed prespecifying the exact solution $\varphi(\boldsymbol{x})$ of the differential problem (2.6), the relative errors on the solution *err* for the cases in Table 2.5 referred to $\varphi1$ and $\varphi8$ become $8.64(-5)$ and $4.68(-5)$ respectively.

Let us solve the weakly nonlinear system (2.10) where $\sigma(\varphi) = \sigma1$ with the simplified Newton–AM method (with 20 iterations for the inner AM procedure)[3] with the stopping criterium

$$\|\boldsymbol{F}(\boldsymbol{w}^{(k+1)})\| \equiv \|\boldsymbol{F}_0(\boldsymbol{w}^{(k+1)})\| \le \tilde{\tau}_a + \tilde{\tau}_r\|\boldsymbol{F}(\boldsymbol{w}^{(0)})\| \tag{2.53}$$

where $\tilde{\tau}_a$ and $\tilde{\tau}_r$ are prefixed absolute and relative error tolerances equal to $0.5 \cdot 10^{-11}$ in order to have, at the final simplified Newton iteration, a residual $\|\boldsymbol{F}(\boldsymbol{w}^{(k+1)})\|$ approximately equal to $9.99 \cdot 10^{-6}$ as in Table 2.7. In this case we obtain 3080 simplified Newton method's iterations with the relative error *err* $= 4.06 \cdot 10^{-9}$, the residual *res* $= 1.05 \cdot 10^{-5}$ and the initial residual *res0* $= 2.1 \cdot 10^6$.

From the numerical experiments the following conclusions can be drawn:

- from tables 2.1–2.3, we can observe that when the values of the function $\sigma(\varphi)$ increase ($\sigma9$ has larger values than $\sigma7bis$ and $\sigma8$ for $\varphi \in [0, 1]$), then the total number of the simplified Newton iterations increases;

- from tables 2.1–2.4, we observe that when the values of the function $g(\varphi)$ are rapidly increasing for $0 \le \varphi \le 1$ 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;

---

[3]I.e., we consider only one iteration of the LDFI procedure.

- from tables 2.1–2.3, we remark that there is no an appreciable reduction of the total number of the simplified Newton iterations when we solve the weakly nonlinear system with a greater precision than that imposed on the LDFI–procedure. Thus, the first and the second iteration levels should have approximately the same precision on the stopping criterium;

- from tables 2.9–2.11, we remark that the LDFI–procedure combined with the simplified Newton–AM method gives better results (in terms of total number of Newton iterations and of the number of LDFI–procedure iterations) when the coefficients of $\tilde{\boldsymbol{v}}$ increase, i.e., when the *deviation from asymmetry*[4] of the matrix $C_\nu$ increases. This is a peculiar feature of the Arithmetic Mean linear solver ([37], [14]) especially when it is implemented as inner solver in a two iteration levels procedure, such as the Newton–AM method ([11], [13]).

In the case of 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.

$N = 256;\ \sigma(\varphi) = \sigma 7_1;\ \varphi(\boldsymbol{x}) = \varphi 8;\ \alpha(x, y) = 0;\ \tilde{v}_1 = \tilde{v}_2 = 10;\ j_k = 20;$

| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| | | $\underline{\varepsilon} = 10^{-5}$ | | | |
| $g0$ | 35 (345) | 1.35(-9) | 8.08(-6) | 8.07(5) | 1.53(-8) |
| $g1_2$ | 34 (201) | 7.02(-10) | 8.91(-6) | 8.08(5) | 8.58(-8) |
| $g1_3$ | 35 (349) | 1.49(-9) | 8.43(-6) | 8.07(5) | 1.59(-8) |
| $g2_2$ | 34 (300) | 9.47(-10) | 9.38(-6) | 8.09(5) | 1.07(-7) |
| $g5_2$ | 35 (327) | 1.37(-9) | 8.58(-6) | 8.07(5) | 1.60(-8) |
| $g5_3$ | 34 (59) | 1.14(-10) | 7.03(-6) | 8.18(5) | 1.72(-8) |
| $g5_4$ | 34 (194) | 8.51(-10) | 9.58(-6) | 8.08(5) | 8.18(-8) |

| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| | | $\underline{\varepsilon} = 10^{-3}$ | | | |
| $g0$ | 85 (345) | 1.57(-9) | 9.40(-6) | 8.07(5) | 1.73(-8) |
| $g1_2$ | 63 (201) | 7.65(-10) | 9.43(-6) | 8.08(5) | 1.38(-8) |
| $g1_3$ | 85 (349) | 1.66(-9) | 9.85(-6) | 8.07(5) | 1.82(-8) |
| $g2_2$ | 80 (301) | 9.75(-10) | 9.65(-6) | 8.09(5) | 1.17(-8) |
| $g5_2$ | 82 (327) | 1.59(-9) | 9.97(-6) | 8.07(5) | 1.84(-8) |
| $g5_3$ | 40 (59) | 1.15(-10) | 7.11(-6) | 8.18(5) | 7.02(-9) |
| $g5_4$ | 62 (195) | 8.10(-10) | 8.71(-6) | 8.08(5) | 1.53(-8) |

| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| | | $\underline{\varepsilon} = 10^{-2}$ | | | |
| $g0$ | 108 (345) | 1.63(-9) | 9.77(-6) | 8.07(5) | 1.81(-8) |
| $g1_2$ | 76 (201) | 7.90(-10) | 9.73(-6) | 8.08(5) | 1.43(-8) |
| $g1_3$ | 109 (350) | 1.59(-9) | 9.43(-6) | 8.07(5) | 1.75(-8) |
| $g2_2$ | 102 (302) | 9.21(-10) | 9.15(-6) | 8.09(5) | 1.11(-8) |
| $g5_2$ | 104 (328) | 1.51(-9) | 9.49(-6) | 8.07(5) | 1.76(-8) |
| $g5_3$ | 42 (59) | 1.16(-10) | 7.14(-6) | 8.18(5) | 7.04(-9) |
| $g5_4$ | 75 (195) | 8.38(-10) | 9.01(-6) | 8.08(5) | 1.60(-8) |

Table 2.1: Results for different functions $g(\varphi)$ and different value of $\underline{\varepsilon}$: $\sigma(\varphi) = \sigma 7_1$.

---

[4]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 $A$ ([11]).

$N = 256$; $\sigma(\varphi) = \sigma8$; $\varphi(\boldsymbol{x}) = \varphi8$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 10$; $j_k = 20$;

$\underline{\varepsilon} = 10^{-5}$

| $g(\varphi)$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| $g0$ | 34 (304) | 1.76(-9) | 9.02(-6) | 8.07(5) | 1.80(-7) |
| $g1_2$ | 34 (174) | 7.60(-10) | 9.40(-6) | 8.08(5) | 8.10(-8) |
| $g1_3$ | 34 (308) | 1.80(-9) | 9.09(-6) | 8.07(5) | 1.82(-7) |
| $g2_2$ | 34 (266) | 1.11(-9) | 9.49(-6) | 8.09(5) | 9.85(-8) |
| $g5_2$ | 34 (288) | 1.63(-9) | 8.85(-6) | 8.07(5) | 2.11(-7) |
| $g5_3$ | 34 (51) | 9.48(-11) | 6.99(-6) | 8.18(5) | 1.75(-8) |
| $g5_4$ | 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 |
|---|---|---|---|---|---|
| $g0$ | 77 (305) | 1.87(-9) | 9.70(-6) | 8.07(5) | 2.39(-8) |
| $g1_2$ | 60 (175) | 7.20(-10) | 9.02(-6) | 8.08(5) | 1.46(-8) |
| $g1_3$ | 77 (309) | 1.93(-9) | 9.83(-6) | 8.07(5) | 2.44(-8) |
| $g2_2$ | 74 (268) | 1.06(-9) | 9.10(-6) | 8.09(5) | 1.43(-8) |
| $g5_2$ | 75 (289) | 1.70(-9) | 9.35(-6) | 8.07(5) | 2.26(-8) |
| $g5_3$ | 38 (51) | 9.63(-11) | 7.13(-6) | 8.18(5) | 6.94(-9) |
| $g5_4$ | 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 |
|---|---|---|---|---|---|
| $g0$ | 97 (306) | 1.77(-9) | 9.19(-6) | 8.07(5) | 2.29(-8) |
| $g1_2$ | 71 (175) | 7.45(-10) | 9.34(-6) | 8.08(5) | 1.52(-8) |
| $g1_3$ | 98 (310) | 1.82(-9) | 9.30(-6) | 8.07(5) | 2.33(-8) |
| $g2_2$ | 92 (268) | 1.09(-9) | 9.45(-6) | 8.09(5) | 1.51(-8) |
| $g5_2$ | 94 (289) | 1.77(-9) | 9.73(-6) | 8.07(5) | 2.38(-8) |
| $g5_3$ | 40 (51) | 9.72(-11) | 7.20(-6) | 8.18(5) | 7.01(-9) |
| $g5_4$ | 69 (168) | 7.64(-10) | 8.79(-6) | 8.08(5) | 1.64(-8) |

Table 2.2: Results for different functions $g(\varphi)$ and different value of $\underline{\varepsilon}$: $\sigma(\varphi) = \sigma8$.

$N = 256;\ \sigma(\varphi) = \sigma 9;\ \varphi(\boldsymbol{x}) = \varphi 8;\ \alpha(x,y) = 0;\ \tilde{v}_1 = \tilde{v}_2 = 10;\ j_k = 20;$

$\underline{\varepsilon = 10^{-5}}$

| $g(\varphi)$ | $it(ktot)$ | $err$ | $res$ | $res0$ | $diff$ |
|---|---|---|---|---|---|
| $g0$ | 47 (3502) | 9.94(-10) | 9.93(-6) | 7.47(7) | 1.09(-9) |
| $g1_2$ | 44 (1578) | 4.77(-10) | 9.97(-6) | 7.47(7) | 1.03(-9) |
| $g1_3$ | 47 (3573) | 1.01(-9) | 9.94(-6) | 7.47(7) | 1.09(-9) |
| $g2_2$ | 45 (2660) | 8.25(-10) | 9.95(-6) | 7.47(7) | 1.12(-9) |
| $g5_2$ | 46 (3214) | 9.27(-10) | 9.92(-6) | 7.47(7) | 1.09(-9) |
| $g5_3$ | 43 (579) | 1.20(-10) | 9.59(-6) | 7.48(7) | 6.94(-10) |
| $g5_4$ | 44 (1550) | 4.56(-10) | 9.95(-6) | 7.47(7) | 1.01(-9) |

$\underline{\varepsilon = 10^{-3}}$

| $g(\varphi)$ | $it(ktot)$ | $err$ | $res$ | $res0$ | $diff$ |
|---|---|---|---|---|---|
| $g0$ | 607 (3274) | 1.01(-9) | 9.97(-6) | 7.47(7) | 1.05(-9) |
| $g1_2$ | 326 (1526) | 4.68(-10) | 9.87(-6) | 7.47(7) | 9.60(-10) |
| $g1_3$ | 615 (3337) | 1.03(-9) | 9.97(-6) | 7.47(7) | 1.05(-9) |
| $g2_2$ | 499 (2530) | 8.28(-10) | 9.92(-6) | 7.47(7) | 1.06(-9) |
| $g5_2$ | 570 (3019) | 9.44(-10) | 9.98(-6) | 7.47(7) | 1.04(-9) |
| $g5_3$ | 147 (569) | 1.15(-10) | 9.64(-6) | 7.48(7) | 6.29(-10) |
| $g5_4$ | 320 (1498) | 4.94(-10) | 9.92(-6) | 7.47(7) | 9.41(-10) |

$\underline{\varepsilon = 10^{-2}}$

| $g(\varphi)$ | $it(ktot)$ | $err$ | $res$ | $res0$ | $diff$ |
|---|---|---|---|---|---|
| $g0$ | 890 (3162) | 1.01(-9) | 9.99(-6) | 7.47(7) | 1.05(-9) |
| $g1_2$ | 468 (1500) | 4.69(-10) | 9.87(-6) | 7.47(7) | 9.62(-10) |
| $g1_3$ | 904 (3221) | 1.03(-9) | 9.96(-6) | 7.47(7) | 1.05(-9) |
| $g2_2$ | 728 (2466) | 8.28(-10) | 9.90(-6) | 7.47(7) | 1.05(-9) |
| $g5_2$ | 836 (2924) | 9.41(-10) | 9.94(-6) | 7.47(7) | 1.04(-9) |
| $g5_3$ | 199 (564) | 1.15(-10) | 9.65(-6) | 7.48(7) | 6.29(-10) |
| $g5_4$ | 459 (1472) | 4.52(-10) | 9.93(-6) | 7.47(7) | 9.44(-10) |

Table 2.3: Results for different functions $g(\varphi)$ and different value of $\underline{\varepsilon}$: $\sigma(\varphi) = \sigma 9$.

$\sigma(\varphi) = \sigma 7_2;\ g(\varphi) = g5_2;\ \varphi(\boldsymbol{x}) = \varphi 8;\ \tilde{v}_1 = \tilde{v}_2 = 100;\ j_k = 20;$

| $\alpha(\boldsymbol{x})$ | $c$ | $it(ktot)$ | $err$ | $res$ | $res0$ | $diff$ |
|---|---|---|---|---|---|---|
| 0 | | 33 (62) | 1.21(-11) | 3.26(-6) | 1.84(6) | 3.38(-9) |
| $\alpha 1$ | 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) |
| $\alpha 2$ | 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) |

Table 2.4: Results for different functions $\alpha(\boldsymbol{x})$.

$N = 256$; $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_1$; $\alpha(x,y) = 0$; $\tilde{v}_1 00 = \tilde{v}_2 = 100$; $j_k = 20$;
Homogeneous boundary conditions;

| $\varphi(\boldsymbol{x})$ | it(ktot) | err | res | res0 |
|---|---|---|---|---|
| $\varphi 1$ | 21 (46) | 4.18(-11) | 1.59(-6) | 1.84(6) |
| $\varphi 21$ | 13 (26) | 8.22(-12) | 5.09(-7) | 5.22(5) |
| $\varphi 22$ | 21 (34) | 3.09(-11) | 7.28(-6) | 5.26(5) |
| $\varphi 23$ | 26 (37) | 4.14(-12) | 9.05(-6) | 3.17(5) |
| $\varphi 31$ | 15 (28) | 2.57(-11) | 3.17(-6) | 5.23(5) |
| $\varphi 32$ | 19 (32) | 4.47(-11) | 8.48(-6) | 5.22(5) |
| $\varphi 8$ | 33 (62) | 1.29(-11) | 3.48(-6) | 1.84(6) |

$N = 256$; $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_1$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 1$; $j_k = 20$
Nonhomogeneous boundary conditions;

| $\varphi(\boldsymbol{x})$ | it(ktot) | err | res | res0 |
|---|---|---|---|---|
| $\varphi 4$ | 36 (3016) | 1.91(-8) | 9.96(-6) | 1.73(6) |
| $\varphi 5$ | 34 (2234) | 1.49(-8) | 9.94(-6) | 4.19(5) |
| $\varphi 6$ | 29 (184) | 6.50(-8) | 8.21(-6) | 1.17(4) |
| $\varphi 7$ | 27 (178) | 3.42(-7) | 7.15(-6) | 2.53(3) |

Table 2.5: Results for different functions $\varphi(\boldsymbol{x})$.

$N = 256$; $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_2$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 100$; $j_k = 20$;

| $\varphi(\boldsymbol{x})$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| $\varphi 1$ | 21 (46) | 4.14(-11) | 1.58(-6) | 1.84(6) | 2.89(-9) |
| $\varphi 2_1$ | 13 (26) | 5.32(-12) | 3.95(-7) | 5.22(5) | 4.34(-8) |
| $\varphi 2_2$ | 20 (33) | 3.50(-11) | 7.14(-6) | 5.26(5) | 2.54(-8) |
| $\varphi 2_3$ | 26 (37) | 3.57(-12) | 7.88(-6) | 3.17(5) | 1.22(-8) |
| $\varphi 3_1$ | 15 (28) | 2.45(-11) | 3.03(-6) | 5.23(5) | 5.83(-8) |
| $\varphi 3_2$ | 19 (32) | 5.23(-12) | 1.28(-6) | 5.22(5) | 9.47(-9) |
| $\varphi 8$ | 33 (62) | 1.21(-11) | 3.26(-6) | 1.84(6) | 3.38(-9) |

Table 2.6: Results for different functions $\varphi(\boldsymbol{x})$.

$N = 256$; $g(\varphi) = g5_1$; $\varphi(\boldsymbol{x}) = \varphi 8$; $\alpha(x,y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 1$; $j_k = 20$;

| $\sigma(\varphi)$ | it(ktot) | err | res | res0 |
|---|---|---|---|---|
| $\sigma 1$ | 36 (3089) | 3.84(-9) | 9.99(-6) | 2.12(6) |
| $\sigma 2$ | 36 (1760) | 7.33(-9) | 8.09(-6) | 7.62(5) |
| $\sigma 3$ | 35 (1856) | 5.66(-9) | 9.31(-6) | 1.65(6) |
| $\sigma 4$ | 39 (3282) | 6.25(-9) | 9.97(-6) | 1.62(6) |
| $\sigma 5$ | 39 (2705) | 3.62(-9) | 8.47(-6) | 1.51(6) |
| $\sigma 6$ | 39 (2962) | 3.98(-9) | 8.29(-6) | 1.51(6) |
| $\sigma 7_1$ | 36 (1761) | 3.53(-9) | 7.18(-6) | 7.82(5) |
| $\sigma 7_2$ | 35 (2785) | 3.72(-9) | 8.51(-6) | 1.67(6) |
| $\sigma 8$ | 34 (1458) | 1.11(-8) | 2.32(-6) | 7.82(5) |
| $\sigma 9$ | 47 (4846) | 1.55(-9) | 9.99(-6) | 7.48(7) |

Table 2.7: Results for different functions $\sigma(\varphi)$.

$\sigma(\varphi) = \sigma 7_2$; $\varphi(\boldsymbol{x}) = \varphi 8$; $\alpha(x, y) = 0$; $\tilde{v}_1 = \tilde{v}_2 = 0$; $j_k = 20$;

| | | $g(\varphi) = g5_1$ | | | |
|---|---|---|---|---|---|
| $N$ | it(ktot) | err | res | res0 | diff |
| 32 | 31 (38) | 4.27(-9) | 4.80(-6) | 4.63(3) | 1.27(-7) |
| 64 | 35 (333) | 2.96(-8) | 9.17(-6) | 2.51(4) | 1.64(-7) |
| 128 | 33 (425) | 1.03(-8) | 9.00(-6) | 1.39(5) | 5.66(-8) |
| 256 | 35 (1866) | 4.48(-9) | 9.52(-6) | 7.79(5) | 2.00(-8) |

| | | $g(\varphi) = g5_2$ | | | |
|---|---|---|---|---|---|
| $N$ | 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) |

Table 2.8: Results for different values of $N$.

$N = 256$, $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_1$; $\varphi(\boldsymbol{x}) = \varphi 2_3$; $\alpha(x, y) = 0$; $j_k = 20$;

| $\tilde{\boldsymbol{v}}$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| $(0, 0)^T$ | 34 (2266) | 8.12(-9) | 9.83(-6) | 2.84(5) | 1.19(-8) |
| $(1, 1)^T$ | 34 (2093) | 7.39(-9) | 9.84(-6) | 2.84(5) | 1.19(-8) |
| $(10, 10)^T$ | 33 (384) | 7.56(-10) | 9.94(-6) | 2.87(5) | 5.39(-8) |
| $(50, 50)^T$ | 33 (73) | 4.45(-11) | 6.55(-6) | 2.98(5) | 1.14(-8) |
| $(100, 0)^T$ | 28 (42) | 7.77(-12) | 5.38(-6) | 3.01(5) | 8.65(-9) |
| $(0, 100)^T$ | 28 (43) | 2.03(-12) | 2.09(-6) | 3.01(5) | 3.19(-9) |
| $(100, 50)^T$ | 28 (41) | 2.04(-11) | 9.28(-6) | 3.08(5) | 1.60(-8) |
| $(50, 100)^T$ | 27 (42) | 8.36(-12) | 7.94(-6) | 3.08(5) | 1.28(-8) |
| $(100, 100)^T$ | 26 (37) | 4.14(-12) | 9.05(-6) | 3.17(5) | 1.41(-8) |
| $(150, 150)^T$ | 22 (28) | 3.41(-12) | 3.64(-6) | 3.40(5) | 6.71(-9) |

Table 2.9: Results for different values of $\tilde{\boldsymbol{v}}$.

$N = 256$, $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_2$; $\varphi(\boldsymbol{x}) = \varphi 8$; $\alpha(x, y) = 0$; $j_k = 20$;

| $\tilde{\boldsymbol{v}}$ | it(ktot) | err | res | res0 | diff |
|---|---|---|---|---|---|
| $(0, 0)^T$ | 35 (2318) | 3.14(-9) | 8.67(-6) | 1.67(6) | 3.78(-7) |
| $(1, 1)^T$ | 35 (2269) | 3.10(-9) | 8.75(-6) | 1.67(6) | 3.65(-7) |
| $(5, 5)^T$ | 35 (1552) | 2.15(-9) | 9.20(-6) | 1.68(6) | 2.60(-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) |
| $(100, 0)^T$ | 35 (79) | 3.28(-11) | 5.65(-6) | 1.76(6) | 4.24(-9) |
| $(0, 100)^T$ | 35 (78) | 3.34(-11) | 5.94(-6) | 1.76(6) | 4.46(-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 2.10: Results for different values of $\tilde{\boldsymbol{v}}$.

$N = 256$, $\sigma(\varphi) = \sigma 7_2$; $g(\varphi) = g5_2$; $\varphi(\boldsymbol{x}) = \varphi 8$; $\alpha(x, y) = 0$;

| $\tilde{\boldsymbol{v}}$ | $j_k = 1$ | $j_k = 5$ | $j_k = 10$ | $j_k = 20$ | $j_k = k+1$ | $j_k = \nu+1$ |
|---|---|---|---|---|---|---|
| $(0,0)^T$ | 35 (46353) | 35 (9271) | 35 (4636) | 35 (2318) | 35 (1667) | 35 (2595) |
| | 46353 | 46355 | 46360 | 46360 | 46411 | 46348 |
| | 3.15(-9) | 3.15(-9) | 3.14(-9) | 3.14(-9) | 3.08(-9) | 3.16(-9) |
| $(5,5)^T$ | 35 (31033) | 35(6207) | 35 (3104) | 35 (1552) | 35 (1372) | 35 (1870) |
| | 31033 | 31035 | 31040 | 31040 | 31031 | 31037 |
| | 2.16(-9) | 2.15(-9) | 2.15(-9) | 2.15(-9) | 2.16(-9) | 2.15(-9) |
| $(10,10)^T$ | 35 (16978) | 35 (3396) | 35 (1698) | 35 (849) | 35 (1018) | 35 (1172) |
| | 16978 | 16980 | 16980 | 16980 | 17003 | 16988 |
| | 1.19(-9) | 1.18(-9) | 1.18(-9) | 1.18(-9) | 1.15(-9) | 1.17(-9) |
| $(50,50)^T$ | 36 (2533) | 35 (507) | 35 (254) | 35 (126) | 35 (369) | 35 (326) |
| | 2533 | 2535 | 2540 | 2520 | 2553 | 2543 |
| | 1.13(-10) | 1.09(-10) | 1.08(-10) | 1.13(-10) | 1.07(-10) | 9.07(-11) |
| $(100,50)^T$ | 36 (1465) | 36 (294) | 36 (148) | 35 (74) | 35 (272) | 32 (235) |
| | 1465 | 1470 | 1480 | 1480 | 1471 | 1486 |
| | 5.13(-11) | 4.50(-11) | 3.57(-11) | 2.92(-11) | 4.57(-11) | 2.21(-11) |
| $(50,100)^T$ | 36 (1463) | 35 (293) | 35 (147) | 34 (74) | 35 (271) | 32 (237) |
| | 1463 | 1465 | 1470 | 1480 | 1468 | 1476 |
| | 5.06(-11) | 4.88(-11) | 4.72(-11) | 4.24(-11) | 4.25(-11) | 3.08(-11) |
| $(100,100)^T$ | 36 (1213) | 35 (243) | 36 (122) | 33 (62) | 35 (239) | 30 (218) |
| | 1213 | 1215 | 1220 | 1240 | 1218 | 1226 |
| | 3.78(-11) | 3.46(-11) | 2.31(-11) | 1.21(-11) | 3.59(-11) | 1.78(-11) |
| $(150,150)^T$ | 36 (808) | 36 (162) | 33 (82) | 23 (41) | 35 (186) | 26 (184) |
| | 808 | 810 | 820 | 820 | 807 | 811 |
| | 1.99(-11) | 1.35(-11) | 1.04(-11) | 1.43(-11) | 1.75(-11) | 1.95(-11) |

Table 2.11: In a column are: *it(ktot)*, *jtot* and *err* for different values of $j_k$.

# Bibliography

[1] Aris R.: The Mathematical Theory of Diffusion and Reaction in Permeable Catalysts, vol. I, II, Clarendon Press, Oxford, 1975.

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

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

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

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

[6] Buffoni G., Griffa A., Li Z., de Mottoni P.: *Spatially distributed communities: the resource–consumer system*, Journal of Mathematical Biology, 33 (1995), 723–743.

[7] Chawla M.M., Al–Zanaidi M.A.: *An extended trapezoidal formula for the diffusion equation in two space dimensions*, Computers Mathematics with Applications, 42 (2001), 157–168.

[8] Courant R., Friedrichs K.O., Lewy H.: *Über die partiellen Differenzengleichungen der mathematischen Physik*, Mathematische Annalen, 100 (1928), 32–74.

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

[10] Galligani E.: *A two-stage iterative method for solving a weakly nonlinear parametrized system*, International Journal of Computer Mathematics, 79 (2002), 1211–1224.

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

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

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

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

[15] 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.

[16] 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.

[17] Galligani E., Ruggiero V.: *Analysis of splitting parallel methods for solving block tridiagonal linear systems*, in: Proceedings SMS TPE'94 (V.P. Ivannikov, V.A. Serebriakov eds.), Russian Academy of Sciences, Moscow, 1994, 406–416.

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

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

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

[21] Henrici P.: Discrete Variable Methods in Ordinary Differential Equations, John Wiley & Sons Inc., New York, 1962.

[22] Ito K., Kunisch K.: *Augmented Lagrangian–SQP methods for nonlinear optimal control problems of tracking–type*, SIAM Journal on Optimization, 6 (1996), 29–43.

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

[24] Kelley C.T.: Iterative Methods for Linear and Nonlinear Equations, SIAM, Philadelphia, 1995.

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

[26] Kunisch K., Volkwein S.: *Augmented Lagrangian–SQP techniques and their approximations*, Contemporary Mathematics, 209 (1997), 147–159.

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

[28] 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.

[29] 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.

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

[31] 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.

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

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

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

[35] 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.

[36] Pohozaev S.T.: *The Dirichlet problem for the equation $\Delta u = u^2$*, Soviet Mathematics, 1 (1960), 1143–1146.

[37] Ruggiero V., Galligani E.: *An iterative method for large sparse linear systems on a vector computer*, Computers and Mathematics with Applications, 20 (1990), 25–28.

[38] Sawami H., Niki H.: *On iterative methods for solving a semi–linear eigenvalue problem*, International Journal of Computer Mathematics, 62 (1996), 271–284.

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

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

[41] 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.

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