

# DISCRETE MAXIMUM PRINCIPLES FOR THE COURANT FINITE ELEMENT SOLUTION OF SOME NONLINEAR ELLIPTIC PROBLEMS

Menghis T. Bahlibi (Budapest, Hungary)

Communicated by László Szili

(Received September 30, 2023; accepted January 7, 2024)

**Abstract.** The discrete maximum principle (DMP) is an important measure of the qualitative reliability of the numerical solutions for elliptic PDE models. We are motivated by known results on the DMP and nonnegativity preservation of finite element (FE) solutions under the condition that sufficiently small enough mesh size  $h$  is used. We extend the above results by explicitly looking for how much the mesh size  $h$  should be small to guarantee the qualitative properties. We determine a threshold mesh size  $h_0$  in terms of the angle condition to ensure the validity of DMPs by Courant FEM.

## 1. Introduction

The maximum principle forms an important qualitative property of second-order elliptic equations [11], therefore its discrete analogs, the so-called DMPs have been studied by many researchers [1, 2, 3, 6, 12]. The DMP is an important measure of the qualitative reliability of the numerical scheme, otherwise one could get unphysical numerical solutions like negative concentrations, etc. Typical maximum principles arise either in the form

$$\max_{\bar{\Omega}} u = \max_{\Gamma_D} u$$

---

*Key words and phrases:* Nonlinear elliptic problem, discrete maximum principle, finite element method, angle condition on the mesh.

*2010 Mathematics Subject Classification:* 65N30, 35J60.

<https://doi.org/10.71352/ac.57.055>

(that is, the solution  $u$  attains its maximum on the boundary) or in the form of

$$\max_{\bar{\Omega}} u \leq \max\{0, \max_{\Gamma_D} u\}$$

(that is, the solution  $u$  can attain a nonnegative maximum only on the boundary). Analogous minimum principles are defined by reversing signs. A physically important special case is nonnegativity preservation.

We are motivated by the articles [6, 7, 8] which deal with this topic and use a sufficiently small mesh size to establish the qualitative properties, in particular, DMP and nonnegativity. Compared to this, we explicitly look for how much the step size  $h$  must be small to guarantee the qualitative properties. We determine a threshold mesh size  $h_0$  in terms of the angle condition to ensure the validity of DMPs by Courant FEM.

We formulate and prove the corresponding DMPs. Finally, an example of a real-life problem, where the preservation of maximum principles plays an important role, is presented.

## 2. The problem and its discretization

Let us begin with the following nonlinear elliptic model:

$$(2.1) \quad \begin{cases} -\operatorname{div} \left( b(x, u, \nabla u) \nabla u \right) + r(x, u, \nabla u) u = f(x) & \text{in } \Omega, \\ b(x, u, \nabla u) \frac{\partial u}{\partial \nu} = \gamma(x) & \text{on } \Gamma_N, \\ u = g(x) & \text{on } \Gamma_D, \end{cases}$$

where  $\Omega$  is a bounded domain in  $\mathbb{R}^2$  under the assumptions below:

- (a)  $\Omega$  has a piecewise smooth and Lipschitz continuous boundary  $\partial\Omega$ ;  $\Gamma_N$ ,  $\Gamma_D \subset \partial\Omega$  are measurable open sets, such that  $\Gamma_N \cap \Gamma_D = \emptyset$ ,  $\bar{\Gamma}_N \cup \bar{\Gamma}_D = \partial\Omega$  and  $\operatorname{meas}(\Gamma_D) > 0$ .
- (b) The scalar functions  $b : \bar{\Omega} \times \mathbb{R} \times \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $r : \bar{\Omega} \times \mathbb{R} \times \mathbb{R}^2 \rightarrow \mathbb{R}$  are continuous. Further,  $f \in L^2(\Omega)$ ,  $\gamma \in L^2(\Gamma_N)$  and  $g = g^*|_{\Gamma_D}$  with  $g^* \in H^1(\Omega)$ .
- (c) The functions  $b$  and  $r$  are bounded such that

$$(2.2) \quad 0 < \mu_0 \leq b(x, \xi, \eta) \leq \mu_1, \quad 0 \leq r(x, \xi, \eta) \leq \beta \quad \forall (x, \xi, \eta) \in \bar{\Omega} \times \mathbb{R} \times \mathbb{R}^2,$$

where  $\mu_0$ ,  $\mu_1$  and  $\beta$  are positive constants. The weak formulation of the problem (2.1) and its unique weak solution  $u \in H^1(\Omega)$  are defined as follows:

$$(2.3) \quad u = g \quad \text{on } \Gamma_D \text{ in trace sense} \quad \text{and}$$

$$\begin{aligned}
(2.4) \quad & \int_{\Omega} [b(x, u, \nabla u) \nabla u \cdot \nabla v + r(x, u, \nabla u)u] dx = \\
& = \int_{\Omega} f v dx + \int_{\Gamma_N} \gamma v d\sigma \quad \forall v \in H_D^1(\Omega).
\end{aligned}$$

To find the finite element solution, we solve the following problem (which is the counterpart of (2.3)–(2.4) in  $V_h$ ): Find  $u_h \in V_h$  such that

$$u_h = g_h \text{ on } \Gamma_D \quad \text{and}$$

$$\begin{aligned}
(2.5) \quad & \int_{\Omega} [b(x, u_h, \nabla u_h) \nabla u_h \cdot \nabla v_h + r(x, u_h, \nabla u_h)u_h v_h] dx = \\
& = \int_{\Omega} f_h v_h dx + \int_{\Gamma_N} \gamma_h v_h d\sigma
\end{aligned}$$

$\forall v_h \in V_h^0$ . Let the vector  $\bar{\mathbf{c}} = (c_1, \dots, c_{n+m})^T$  contain the values of the finite element solution  $u_h$  at all the nodal points i.e.  $c_i = u_h(P_i)$  and  $u_h = \sum_{i=1}^{n+m} c_i \phi_i$ , where  $\phi_1, \dots, \phi_n$  are the interior basis functions and  $\phi_{n+1}, \dots, \phi_{n+m}$  are the boundary basis functions. Furthermore,  $\bar{\mathbf{b}} = (b_1, \dots, b_n, g_1, \dots, g_m)^T$ . Here, a nonlinear algebraic system of equations is obtained:

$$(2.6) \quad \bar{\mathbf{A}}(\bar{\mathbf{c}})\bar{\mathbf{c}} = \bar{\mathbf{b}},$$

and the structure of the matrix in (2.6) is

$$(2.7) \quad \bar{\mathbf{A}}(\bar{\mathbf{c}}) = \begin{pmatrix} \mathbf{A}(\bar{\mathbf{c}}) & \tilde{\mathbf{A}}(\bar{\mathbf{c}}) \\ \mathbf{0} & \mathbf{I} \end{pmatrix},$$

where  $\mathbf{I}$  is an  $m \times m$  identity matrix and  $\mathbf{0}$  is a  $m \times n$  zero matrix and  $\bar{\mathbf{A}}(\bar{\mathbf{c}})$  is  $(n+m)$  by  $(n+m)$  matrix.

### 3. DMPs for linear systems

If equation (2.5) is in particular linear (when  $b$  and  $r$  are independent of  $u$ ), then the algebraic system of the equations and the structure of the matrix are respectively

$$(3.1) \quad \bar{\mathbf{A}}\bar{\mathbf{c}} = \bar{\mathbf{b}} \quad \text{and} \quad \bar{\mathbf{A}} = \begin{pmatrix} \mathbf{A} & \tilde{\mathbf{A}} \\ \mathbf{0} & \mathbf{I} \end{pmatrix},$$

where the matrix  $\bar{\mathbf{A}}$  has a dimension of  $(n + m)$  by  $(n + m)$ . The DMPs for such linear systems have been studied, e.g., in [5, 9].

**Definition 3.1.** The matrix  $\bar{\mathbf{A}}$  in (3.1) satisfies

- the *discrete weak maximum principle* (DwMP) if for any vector  $\bar{\mathbf{c}} = (c_1, \dots, c_{n+m})^T \in \mathbb{R}^{n+m}$  satisfying  $(\bar{\mathbf{A}}\bar{\mathbf{c}})_i \leq 0$ ,  $i = 1, \dots, n$ , one has

$$\max_{i=1, \dots, n+m} c_i \leq \max\{0, \max_{i=n+1, \dots, n+m} c_i\};$$

- the *discrete strict weak maximum principle* (DWMP) if for any vector  $\bar{\mathbf{c}} = (c_1, \dots, c_{n+m})^T \in \mathbb{R}^{n+m}$  satisfying  $(\bar{\mathbf{A}}\bar{\mathbf{c}})_i \leq 0$ ,  $i = 1, \dots, n$ , one has

$$\max_{i=1, \dots, n+m} c_i = \max_{i=n+1, \dots, n+m} c_i.$$

**Theorem 3.1.** Let the matrix  $\bar{\mathbf{A}}$  in (3.1) satisfy the following conditions, where  $a_{ij}$  denote the entries of  $\bar{\mathbf{A}}$ :

- (i)  $a_{ij} \leq 0 \quad (\forall i = 1, \dots, n, j = 1, \dots, n + m; i \neq j),$
- (ii)  $\sum_{j=1}^{n+m} a_{ij} \geq 0 \quad (\forall i = 1, \dots, n),$
- (iii)  $\mathbf{A}$  is positive definite.

Then  $\bar{\mathbf{A}}$  possesses the DwMP. If the inequality in condition (ii) is replaced by equality, then  $\bar{\mathbf{A}}$  possesses the DWMP.

This theorem is proved in [6].

## 4. DMPs for nonlinear elliptic problems

Now, we are ready to state theorems related to problem (2.1).

### 4.1. The general result

**Definition 4.1.** The family  $\mathcal{F}$  of triangulations of triangular meshes of a bounded polygonal domain is said to be *uniformly acute* if there exists  $\alpha_0 < \frac{\pi}{2}$  such that  $\alpha_n \leq \alpha_0$  for any angle  $\alpha_n$  in all  $T_k$  in all  $\mathcal{T}_h$ , where  $\mathcal{T}_h \in \mathcal{F}$ .

**Definition 4.2.** The *mesh width*  $h$  is the longest diameter occurring in the triangulation  $\mathcal{T}_h$ . i.e.  $h := \max_{k=1, \dots, M} \text{diam}(T_k)$ , where  $\text{diam}$  refers to diameter.

**Theorem 4.1.** *Let the conditions (a)–(c) hold and the Courant finite element method be used with triangulations satisfying Definition 4.1. Let the mesh size  $h$  satisfy*

$$(4.1) \quad 0 < h \leq h_0 = \left( \frac{12 \cos(\alpha_0) \mu_0}{\beta} \right)^{\frac{1}{2}},$$

where  $\alpha_0$  is the angle that obeys Definition 4.1,  $\mu_0$  and  $\beta$  are positive constants from (2.2). Then the matrix in (2.7) satisfies the following properties, where  $a_{ij}(\bar{\mathbf{c}})$  denotes the entries of  $\bar{\mathbf{A}}(\bar{\mathbf{c}})$ .

- (i)  $a_{ij}(\bar{\mathbf{c}}) \leq 0, \quad i = 1, \dots, n, \quad j = 1, \dots, n+m \quad (i \neq j),$
- (ii)  $\sum_{j=1}^{n+m} a_{ij}(\bar{\mathbf{c}}) \geq 0, \quad i = 1, \dots, n,$
- (iii)  $\mathbf{A}(\bar{\mathbf{c}})$  is positive definite.

**Proof.** Let  $\phi_i$  and  $\phi_j$  be any basis functions. Then the entries of the matrix  $\bar{\mathbf{A}}(\bar{\mathbf{c}})$  are:

$$(4.2) \quad a_{ij}(\bar{\mathbf{c}}) = \int_{\Omega} \left[ b(x, u_h, \nabla u_h) \nabla \phi_i \cdot \nabla \phi_j + r(x, u_h, \nabla u_h) \phi_i \phi_j \right] dx.$$

We now prove the properties (i)–(iii).

(i) Let  $i = 1, \dots, n, \quad j = 1, \dots, n+m$  with  $i \neq j$  and let  $\Omega_{ij}$  denote the interior of  $\text{supp } \phi_i \cap \text{supp } \phi_j$ . If  $\Omega_{ij} = \emptyset$  then  $a_{ij}(\bar{\mathbf{c}}) = 0$ . If  $\Omega_{ij} \neq \emptyset$ , then to determine (4.2) we should find the bound of the following integrals:

$$(4.3) \quad \int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, dx \quad \text{and} \quad \int_{\Omega} \phi_i \phi_j \, dx.$$

From Definition 4.1 we have the maximum angle  $\alpha_0$  such that  $\cos(\alpha_0) := \sigma_0 \geq 0$  which is independent of  $i, j$  and  $h$ . The goal here is to find an upper bound of the stiffness matrix obtained from the first part of (4.3). Now let us begin with the inner product of the basis functions on a given triangle. For any angle  $\alpha_{ij}$ , we have

$$\begin{aligned} \nabla \phi_i \cdot \nabla \phi_j &= |\nabla \phi_i| \cdot |\nabla \phi_j| \cos(180^\circ - \alpha_{ij}) = \\ &= \frac{1}{h_i} \cdot \frac{1}{h_j} (-\cos(\alpha_{ij})) \leq \frac{-\cos(\alpha_{ij})}{h^2} \leq \\ &\leq \frac{-\cos(\alpha_0)}{h^2} \quad \forall h_i, h_j \leq h, \quad \forall \alpha_{ij} \leq \alpha_0. \end{aligned}$$

Therefore,

$$(4.4) \quad \nabla \phi_i \cdot \nabla \phi_j \leq -\frac{\sigma_0}{h^2} < 0.$$

Hence, using Definition 4.2 and (4.4),

$$(4.5) \quad \int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, dx = \int_{\Omega_{ij}} \nabla \phi_i \cdot \nabla \phi_j \, dx \leq -\frac{\sigma_0}{h^2} \text{meas}(\Omega_{ij}).$$

To estimate the mass matrix obtained from the second part of (4.3) for general triangles, we use a reference triangle. If especially  $E$  is the reference triangle that is placed in the coordinate system with vertices  $(0, 0)$ ,  $(h, 0)$ , and  $(0, h)$  then one can calculate

$$(4.6) \quad \int_E \phi_i \phi_j \, dx = \frac{h^2}{24}.$$

Based on the reference triangle, we can calculate the mass matrix for general triangles  $T_k$  using affine mappings from the reference element onto  $T_k$  such that  $L_k : E \rightarrow T_k$ . We also define  $J_k = L'_k$ . If the reference triangle  $E$  is considered with  $h = 1$  in (4.6) and  $T_k$  is a fixed general triangle then

$$(4.7) \quad \int_{T_k} \phi_i \phi_j \, dx = \det(J_k) \int_E \tilde{\phi}_i \tilde{\phi}_j \, dx = \frac{|T_k|}{12}$$

by change of variables and using the fact that  $\det(J_k) = 2|T_k|$ , see [4], where  $|T_k|$  is the area of the triangle, and  $\tilde{\phi}_i$  and  $\tilde{\phi}_j$  are respectively given by  $\phi_i = \phi_i \circ L_k$ ,  $\tilde{\phi}_j = \phi_j \circ L_k$ . Therefore, (4.7) implies

$$(4.8) \quad \int_{\Omega_{ij}} \phi_i \phi_j \, dx = \sum_{T_k \in \Omega_{ij}} \int_{T_k} \phi_i \phi_j \, dx = \frac{1}{12} \text{meas}(\Omega_{ij}),$$

where  $\Omega_{ij} := \text{supp } \phi_i \cap \text{supp } \phi_j$ . Using (2.2), (4.4) (4.5) and (4.8) in (4.2) we have

$$\begin{aligned} a_{ij}(\bar{\mathbf{c}}) &\leq \mu_0 \int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, dx + \beta \int_{\Omega} \phi_i \phi_j \, dx \leq \\ &\leq -\frac{\sigma_0}{h^2} \mu_0 \text{meas}(\Omega_{ij}) + \frac{\beta}{12} \text{meas}(\Omega_{ij}) = \\ &= \text{meas}(\Omega_{ij}) \left( \frac{-\sigma_0}{h^2} \mu_0 + \frac{\beta}{12} \right). \end{aligned}$$

Let

$$(4.9) \quad a_{ij}(h) := \text{meas}(\Omega_{ij}) \left( -\frac{\sigma_0}{h^2} \mu_0 + \frac{\beta}{12} \right),$$

then

$$(4.10) \quad a_{ij}(\bar{\mathbf{c}}) \leq a_{ij}(h).$$

Therefore, the sum of the terms in the bracket in (4.9) tends to  $-\infty$  as  $h \rightarrow 0$ . (The first term goes to  $-\infty$  as  $h \rightarrow 0$  and the second term remains unchanged). This implies  $a_{ij}(h) \leq 0$  if  $h$  is small. The main task here is to find how much  $h$  should be to get the nonpositivity of (4.2). To determine the threshold  $h = h_0$ , the following equation must hold,

$$-\frac{\sigma_0}{h_0^2} \mu_0 + \frac{\beta}{12} = 0.$$

This implies  $h_0 = \left( \frac{12\sigma_0\mu_0}{\beta} \right)^{\frac{1}{2}}$ . In summary, if  $0 < h \leq h_0 = \left( \frac{12\sigma_0\mu_0}{\beta} \right)^{\frac{1}{2}}$ , then  $a_{ij}(\bar{\mathbf{c}}) \leq 0$  from (4.10).

(ii) For any  $i = 1, \dots, n$ ,

$$(4.11) \quad \begin{aligned} \sum_{j=1}^{n+m} a_{ij}(\bar{\mathbf{c}}) &= \int_{\Omega} b(x, u_h, \nabla u_h) \nabla \phi_i \cdot \nabla \left( \sum_{j=1}^{n+m} \phi_j \right) dx + \\ &+ \int_{\Omega} r(x, u_h, \nabla u_h) \phi_i \left( \sum_{j=1}^{n+m} \phi_j \right) dx = \int_{\Omega} r(x, u_h, \nabla u_h) \phi_i dx \geq 0, \end{aligned}$$

using the fact that  $\sum_{j=1}^{n+m} \phi_j \equiv 1$ ,  $0 \leq \phi_i \leq 1$ , and (2.2).

(iii) To verify that  $\bar{\mathbf{A}}(\bar{\mathbf{c}})$  is positive definite, let  $\mathbf{d} \neq 0$  be an arbitrary vector in  $\mathbb{R}^n$ , formed by the coefficients  $d_i$ , and let

$$v_h = \sum_{j=1}^n d_j \phi_j \in V_h.$$

Then  $v_h \neq 0$ . The vector  $\bar{\mathbf{c}} \in \mathbb{R}^{n+m}$  contains the coefficients for  $u_h$  as given in (2.5)–(2.6). Then, using (2.5) and (2.2), we have

$$A(\bar{\mathbf{c}}) \mathbf{d} \cdot \mathbf{d} = \sum_{i,j=1}^n a_{ij}(\bar{\mathbf{c}}) d_i d_j =$$

$$\begin{aligned}
&= \int_{\Omega} \left( b(x, u_h, \nabla u_h) \nabla \left( \sum_{i=1}^n d_i \phi_i \right) \cdot \nabla \left( \sum_{j=1}^n d_j \phi_j \right) + \right. \\
&\quad \left. + r(x, u_h, \nabla u_h) \sum_{i=1}^n d_i \phi_i \sum_{j=1}^n d_j \phi_j \right) dx = \\
&= \int_{\Omega} \left( b(x, u_h, \nabla u_h) |\nabla v_h|^2 + r(x, u_h, \nabla u_h) v_h^2 \right) dx \geq \\
&\geq \mu_0 \int_{\Omega} |\nabla v_h|^2 = \mu_0 |v_h|_1^2 > 0.
\end{aligned}$$

Hence,  $A(\bar{\mathbf{c}})$  is a positive definite matrix.

Altogether, equation (4.9) shows that for small enough  $h$ , we have  $a_{ij}(h) \leq 0$ , but we gave a bound  $h_0$  such that  $a_{ij}(h) \leq 0$  for all  $h \leq h_0$  and all  $i, j$ . ■

**Theorem 4.2.** *Under the conditions of Theorem 4.1 and letting*

$$(4.12) \quad f(x) \leq 0 \quad (x \in \Omega) \quad \text{and} \quad \gamma(x) \leq 0 \quad (x \in \Gamma_N),$$

*we have*

$$(4.13) \quad \max_{\bar{\Omega}} u_h \leq \max\{0, \max_{\Gamma_D} g_h\}.$$

*In particular, if  $\max g_h \geq 0$ , then*

$$(4.14) \quad \max_{\bar{\Omega}} u_h = \max_{\Gamma_D} g_h,$$

*and if  $g_h \leq 0$ , then we have the nonpositivity property*

$$(4.15) \quad \max_{\bar{\Omega}} u_h \leq 0.$$

**Proof.** Let  $\bar{\mathbf{c}} = (c_1, \dots, c_{n+m})^T \in \mathbb{R}^{n+m}$  and  $\bar{\mathbf{b}} = (b_1, \dots, b_n, g_1, \dots, g_m)^T \in \mathbb{R}^{n+m}$  be the vectors that appear in (2.6). Then

$$(4.16) \quad (\bar{b})_i = \int_{\Omega} f \phi_i dx + \int_{\Gamma_N} \gamma \phi_i d\sigma \leq 0 \quad (i = 1, \dots, n)$$

owing to  $f \leq 0, \gamma(x) \leq 0$  and  $0 \leq \phi_i \leq 1$ . Then equation (2.6) and (4.16) imply  $\bar{\mathbf{A}}(\bar{\mathbf{c}})\bar{\mathbf{c}} = \bar{\mathbf{b}} \leq 0$ . Using these arguments and the conditions of Theorem 4.1 one can apply Theorem 3.1 to conclude that the matrix  $\bar{\mathbf{A}}(\bar{\mathbf{c}})$  possesses the DwMP, and as a result, we have

$$(4.17) \quad \max_{i=1, \dots, n+m} c_i \leq \max\{0, \max_{i=n+1, \dots, n+m} c_i\} \leq \max\{0, \max_{i=n+1, \dots, n+m} g_i\},$$



since  $c_i = g_i$  for all  $i = n, \dots, n+m$ . Owing to  $0 \leq \phi_i \leq 1$  and  $\sum_{i=1}^{n+m} \phi_i \equiv 1$ , the solution vectors  $u_h$  and  $g_h$  can be estimated as follows:

$$(4.18) \quad u_h = \sum_{i=1}^{n+m} c_i \phi_i \leq \max_{i=1, \dots, n+m} c_i \sum_{i=1}^{n+m} \phi_i \leq \max_{i=1, \dots, n+m} c_i,$$

$$(4.19) \quad \begin{aligned} g_h &= \sum_{i=n+1}^{n+m} g_i \phi_i \leq \max_{i=n+1, \dots, n+m} g_i \sum_{i=n+1}^{n+m} \phi_i \Rightarrow \\ &\Rightarrow \max_{\Gamma_D} g_h \leq \max_{i=n+1, \dots, n+m} g_i. \end{aligned}$$

In the last statement in fact equality holds because of the following arguments. The piecewise linear basis functions satisfy

$$(4.20) \quad \phi_i(P_j) = \delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j, \end{cases}$$

for proper nodes  $P_1, \dots, P_n \in \Omega$  and  $P_{n+1}, \dots, P_{n+m} \in \Gamma_D$ . Let  $g_k := \max_{i=n+1, \dots, n+m} g_i$ . Then, using (4.20)

$$g_h(P_k) = \sum_{i=n+1}^{n+m} g_i \phi_i(P_k) = g_k,$$

since  $\phi_i(P_k) = 1$  and 0 in the other nodes. That is,  $g_h(P_k) = \max_{i=n+1, \dots, n+m} g_i$ . Hence, we obtained the equality

$$\max_{\Gamma_D} g_h = \max_{i=n+1, \dots, n+m} g_i,$$

which implies

$$(4.21) \quad \max\{0, \max_{\Gamma_D} g_h\} = \max\{0, \max_{i=n+1, \dots, n+m} g_i\}.$$

Altogether, (4.17), (4.18) and (4.21) imply (4.13). The last two statements are simple consequences of (4.13). ■

The corresponding discrete minimum principle for the problem (2.1) can be verified in the same way by reversing signs.

**Theorem 4.3.** *Under the conditions of Theorem 4.1 and letting*

$$(4.22) \quad f(x) \geq 0 \quad (x \in \Omega) \quad \text{and} \quad \gamma(x) \geq 0 \quad (x \in \Gamma_N),$$

we have

$$(4.23) \quad \min_{\Omega} u_h \geq \min\{0, \min_{\Gamma_D} g_h\}.$$

In particular, if  $\min g_h \leq 0$ , then

$$(4.24) \quad \min_{\Omega} u_h = \min_{\Gamma_D} g_h,$$

and, if  $g_h \geq 0$ , then we have nonnegativity property

$$(4.25) \quad \min_{\Omega} u_h \geq 0.$$

**Remark 4.1.** Let us apply a uniformly acute triangulation of a bounded polygonal domain for an elliptic problem using the Courant finite element method. For a fixed mesh, it makes sense to check whether  $h \leq \left(\frac{12 \cos(\alpha_0) \mu_0}{\beta}\right)^{\frac{1}{2}}$  to ensure DMP and nonnegativity preservation. If it is not satisfied then you must try with a finer mesh.

## 4.2. Semilinear problems

Let us consider a special case of problem (2.1) to illustrate the theorem:

$$(4.26) \quad \begin{cases} -\operatorname{div} \left( b(x) \nabla u \right) + q(x, u) = f(x) & \text{in } \Omega, \\ b(x) \frac{\partial u}{\partial \nu} = \gamma(x) & \text{on } \Gamma_N, \\ u = g(x) & \text{on } \Gamma_D, \end{cases}$$

where  $\Omega$  is a bounded polygonal domain in  $\mathbb{R}^2$  and  $q \in C^1(\Omega \times \mathbb{R})$ . We assume that there exists  $\beta > 0$  such that

$$(4.27) \quad 0 \leq \frac{\partial q}{\partial \xi}(x, \xi) \leq \beta \quad \text{and} \quad 0 < \mu_0 \leq b(x) \leq \mu_1, \quad \forall (x, \xi) \in (\Omega \times \mathbb{R}).$$

Let us first define a function  $r$  in terms of  $q$  :

$$(4.28) \quad r(x, \xi) := \begin{cases} \frac{q(x, \xi) - q(x, 0)}{\xi}, & \text{if } \xi \neq 0 \\ \frac{\partial q}{\partial \xi}(x, 0), & \text{if } \xi = 0 \end{cases}$$

then

$$(4.29) \quad r(x, \xi) \xi = q(x, \xi) - q(x, 0), \quad \forall \xi \in \mathbb{R}.$$

We need to show that  $r(x, \xi)$  satisfies condition (2.2). For every  $t \in [0, 1]$  the first part of (4.27) implies  $0 \leq \frac{\partial q}{\partial \xi}(x, t\xi) \leq \beta$ , then by Newton–Leibniz Theorem, we have

$$(4.30) \quad 0 \leq \int_0^1 \frac{\partial q}{\partial \xi}(x, t\xi) dt \leq \int_0^1 \beta dt \Rightarrow 0 \leq \frac{q(x, \xi) - q(x, 0)}{\xi} \leq \beta$$

$$\Rightarrow 0 \leq r \leq \beta \quad .$$

Hence,  $0 \leq r(x, \xi) \leq \beta$  on  $\mathbb{R}$ .

Problem (4.26) can be written as:

$$(4.31) \quad \begin{cases} -\operatorname{div} \left( b(x) \nabla u \right) + q(x, u) - q(x, 0) = f(x) - q(x, 0) & \text{in } \Omega, \\ b(x) \frac{\partial u}{\partial \nu} = \gamma(x) & \text{on } \Gamma_N, \\ u = g(x) & \text{on } \Gamma_D. \end{cases}$$

This implies

$$(4.32) \quad \begin{cases} -\operatorname{div} \left( b(x) \nabla u \right) + r(x, u)u = \tilde{f}(x) & \text{in } \Omega, \\ b(x) \frac{\partial u}{\partial \nu} = \gamma(x) & \text{on } \Gamma_N, \\ u = g(x) & \text{on } \Gamma_D, \end{cases}$$

where  $r(x, u)u = q(x, u) - q(x, 0)$  from (4.29) and  $\tilde{f}(x) = f(x) - q(x, 0)$ . Then we can see that (4.32) is a special case of problem (2.1) from Theorem 4.3.

**Corollary 4.1.** *Let us consider problem (4.32) with assumptions of (4.27) and the conditions of Theorem 4.1 that is Definition 4.1 must hold, and  $h$  must satisfy*

$$(4.33) \quad 0 < h \leq h_0 = \left( \frac{12\sigma_0\mu_0}{\beta} \right)^{\frac{1}{2}}.$$

Then:

- if  $\tilde{f} \leq 0$  and  $\gamma \leq 0$ , then

$$(4.34) \quad \max_{\bar{\Omega}} u_h \leq \max\{0, \max_{\Gamma_D} g_h\}.$$

In particular, if  $\max g_h \geq 0$ , then

$$(4.35) \quad \max_{\bar{\Omega}} u_h = \max_{\Gamma_D} g_h,$$

and if  $g_h \leq 0$ , then we have the nonpositivity property

$$(4.36) \quad \max_{\bar{\Omega}} u_h \leq 0.$$

• if  $\tilde{f} \geq 0$  and  $\gamma \geq 0$ , then we have

$$(4.37) \quad \min_{\bar{\Omega}} u_h \geq \min\{0, \min_{\Gamma_D} g_h\}.$$

In particular, if  $\min g_h \leq 0$ , then

$$(4.38) \quad \min_{\bar{\Omega}} u_h = \min_{\Gamma_D} g_h,$$

and if  $g_h \geq 0$ , then we have the nonnegativity property

$$(4.39) \quad \min_{\bar{\Omega}} u_h \geq 0.$$

This corollary is illustrated by the example stated below.

### 4.3. Example: diffusion-kinetics enzyme problem

A diffusion-kinetics equation governing the steady-state concentration  $u$  of some substrate in an enzyme-catalyzed reaction has the following form, see in [10], where we included mixed boundary conditions:

$$(4.40) \quad \begin{cases} \operatorname{div}(D(x)\nabla u) = q(x, u) & \text{in } \Omega, \\ D \frac{\partial u}{\partial n} = 0 & \text{on } \Gamma_N, \\ u = u_0 & \text{on } \Gamma_D, \end{cases}$$

where  $\Omega$  is a bounded domain in  $\mathbb{R}^2$ ,  $D(x)$  is the positive molecular diffusion coefficient of the substrate in a medium containing some continuous distribution of bacteria,  $q$  is the rate of the enzyme-substrate reaction. In this example, the DMP and nonnegativity of the solution are obtained for a particular case for the reaction rate by Michaelis–Menten theory:

$$(4.41) \quad q(x, \xi) = \frac{\epsilon^{-1}\xi}{\xi + k} \quad \text{for } \xi \geq 0,$$

where  $k > 0$  is the Michaelis constant and  $\epsilon > 0$ . The condition of  $D(x)$  is given by  $0 < \mu_0 \leq D(x) \leq \mu_1$ , where  $\mu_0$  and  $\mu_1$  are positive constants. Further,  $u_0 \geq 0$  and  $\beta = \frac{1}{\epsilon k}$ . First, rewrite the equation in the form (4.32). Then  $\tilde{f}(x) = 0$  since  $q(x, 0) = 0$  and  $r(x, \xi) := \frac{\epsilon^{-1}}{\xi + k}$  and it satisfies (4.30) because  $0 \leq r \leq \frac{1}{\epsilon k}$ . We can also extend  $q(x, \xi)$  from  $\xi \geq 0$  to  $\xi \in \mathbb{R}$ , by the formula  $q(x, -\xi) = -q(x, \xi)$  for  $u \leq 0$  [6].

**Corollary 4.2.** *Let us consider the problem (4.40). If the geometric condition in Theorem 4.1 holds and the mesh  $h$  satisfies*

$$(4.42) \quad h \leq h_0 = \left( 12 \cos(\alpha_0) \mu_0 \epsilon k \right)^{\frac{1}{2}},$$

*then the finite element solution  $u_h$  of (4.40) is bounded by:*

$$\max_{\bar{\Omega}} u_h = \max_{\Gamma_D} u_{0h} \quad \text{and} \quad \min_{\bar{\Omega}} u_h \geq 0.$$

**Proof.** (4.40) is a special case of equation (4.26) for  $g := u_0 \geq 0$ ,  $f(x) = 0$ ,  $\gamma(x) = 0$ , and the mesh size  $h$  when  $\beta = \frac{1}{\epsilon k}$  in Theorem 2 equation (4.1) is:

$$h \leq h_0 = \left( \frac{12 \cos(\alpha_0) \mu_0}{(\epsilon k)^{-1}} \right)^{\frac{1}{2}} = \left( 12 \cos(\alpha_0) \mu_0 \epsilon k \right)^{\frac{1}{2}}.$$

Both (4.34) and (4.37) are true, since  $f(x) = 0$ ,  $\gamma(x) = 0$ . Moreover, (4.35) and (4.39) are satisfied because  $g := u_0 \geq 0$ . Therefore,  $\max_{\bar{\Omega}} u_h = \max_{\Gamma_D} u_{0h}$  and  $\min_{\bar{\Omega}} u_h \geq 0 \quad \forall x \in \Omega$ . ■

Altogether, we are able to guarantee the nonnegativity and determine the maximum of the solution of the PDE without knowing the numerical solution, because the maximum of the solution is attained at the boundary.

**Acknowledgements.** I would like to express my deep gratitude to my research supervisor professor János Karátson for his patience, guidance, enthusiastic encouragement, and useful critiques of my research work. I thank the anonymous referee for the useful comments on the manuscript.

## References

- [1] Brandts, J., S. Korotov and M. Křížek, The discrete maximum principle for linear simplicial finite element approximations of a reaction-diffusion problem, *Linear Algebra Appl.*, **429** (2008), 2344–2357.  
<https://doi.org/10.1016/j.laa.2008.06.011>
- [2] Ciarlet, P.G., Discrete maximum principle for finite-difference operators, *Aequationes Math.*, **4** (1970), 338–352.  
<https://doi.org/10.1007/BF01844166>
- [3] Drăgănescu, A., T.F. Dupont and L.R. Scott, Failure of the discrete maximum principle for an elliptic finite element problem, *Math. Comp.* **74** (2005), no. 249, 1–23.  
<https://doi.org/10.1090/S0025-5718-04-01651-5>

- [4] **Elman, H.C., D.J. Silvester and A.J. Wathen**, *Finite Elements and Fast Iterative Solvers: with Applications in Incompressible Fluid Dynamics*, Numerical Mathematics and Scientific Computation, Oxford University Press, New York, 2005.
- [5] **Faragó, I.**, Matrix and Discrete Maximum Principles, in: *LSSC 2009*, LNCS 5910, 563–570 (2010).
- [6] **Karátson J. and S. Korotov**, Discrete maximum principles for finite element solutions of nonlinear elliptic problems with mixed boundary conditions, *Numer. Math.*, **99** (2005), 669–698  
<https://doi.org/10.1007/s00211-004-0559-0>
- [7] **Karátson J. and S. Korotov**, Discrete maximum principles for FEM solutions of some nonlinear elliptic interface problems, *Int. J. Numer. Anal. Modeling.*, **6**(1) (2008), 1–16.
- [8] **Karátson J., S. Korotov and M. Křížek**, On discrete maximum principles for nonlinear elliptic problems, *Math. Comput. Simul.*, **76** (2007), 99–108.  
<https://doi.org/10.1016/j.matcom.2007.01.011>
- [9] **Karátson J. and S. Korotov**, Some discrete maximum principles arising for nonlinear elliptic finite element problems, *Computers and Mathematics with Applications*, **70** (2015), 2732–2741.  
<https://doi.org/10.1016/j.camwa.2015.06.036>
- [10] **Keller, H.B.**, Elliptic boundary value problems suggested by nonlinear diffusion processes, *Archive for Rational Mechanics and Analysis*, **35** (1969), 363–381.  
<https://doi.org/10.1007/BF00247683>
- [11] **Protter, M.H. and H.F. Weinberger**, *Maximum Principles in Differential Equations*, Springer-Verlag, New York, 1984.
- [12] **Vejchodsky, T.**, The discrete maximum principle for Galerkin solutions of elliptic problems, *Cent. Eur. J. Math.*, **10**(1) (2012), 25–43.  
<https://doi.org/10.2478/s11533-011-0085-0>

**Menghis T. Bahlibi**

Department of Applied Analysis and Computational Mathematics

Eötvös Loránd University

Budapest

Hungary

[menghis.teweldebrhan@gmail.com](mailto:menghis.teweldebrhan@gmail.com)