

Computational Investigations of the Discrete Maximum Principles

Tomáš Vejchodský¹ (vejchod@math.cas.cz)

Antti Hannukainen² (antti.hannukainen@hut.fi)

Sergey Korotov² (sergey.korotov@hut.fi)

¹ Institute of Mathematics, Academy of Sciences
Žitná 25, 115 67 Prague 1
Czech Republic

² Institute of Mathematics
Helsinki University of Technology
P.O. Box 1100, FIN-02015 Espoo
Finland



(Continuous) maximum principle



$$-\operatorname{div}(\mathcal{A}\nabla u) + cu = f \quad \text{in } \Omega, \quad u = g \quad \text{on } \partial\Omega$$

$$\text{MaxP :} \quad f \leq 0 \quad \Rightarrow \quad \max_{\Omega} u \leq \max\{0, \max_{\partial\Omega} u\}$$

$$\text{MinP :} \quad f \geq 0 \quad \Rightarrow \quad \min_{\Omega} u \geq \min\{0, \min_{\partial\Omega} u\}$$

$$\text{ComP :} \quad f \geq 0 \ \& \ g \geq 0 \quad \Rightarrow \quad u \geq 0$$

$$\text{MaxP} \quad \Leftrightarrow \quad \text{MinP} \quad \Leftrightarrow \quad \text{ComP}$$

$$-\operatorname{div}(\mathcal{A}\nabla u) + cu = f \quad \text{in } \Omega, \quad u = g \quad \text{on } \partial\Omega$$

- ▶ Weak formulation: $\bar{u} = u + g$

$$u \in V : \quad a(u, v) = (f, v) - a(g, v) \quad \forall v \in V,$$

where $V = H_0^1(\Omega)$, $g \in H^1(\Omega)$

$$a(u, v) = \int_{\Omega} (\nabla u \cdot \nabla v + cuv) \, dx, \quad (f, v) = \int_{\Omega} fv \, dx$$

- ▶ FEM: $\bar{u}_h = u_h + g_h$, $g_h \approx g$, $g \geq 0 \Leftrightarrow g_h \geq 0$

$$u_h \in V_h : \quad a(u_h, v_h) = (f, v_h) - a(g_h, v_h) \quad \forall v_h \in V_h,$$

where $V_h \subset V$ (continuous, piecewise linear on a mesh \mathcal{T}_h)



$$\bar{u}_h = u_h + g_h, \quad u_h \in V_h : \quad a(u_h, v_h) = (f, v_h) - a(g_h, v_h) \quad \forall v_h \in V_h$$

$$\text{DMaxP :} \quad f \leq 0 \quad \Rightarrow \quad \max_{\Omega} \bar{u}_h \leq \max\{0, \max_{\partial\Omega} \bar{u}_h\}$$

$$\text{DMinP :} \quad f \geq 0 \quad \Rightarrow \quad \min_{\Omega} \bar{u}_h \geq \min\{0, \min_{\partial\Omega} \bar{u}_h\}$$

$$\text{DComP :} \quad f \geq 0 \ \& \ g \geq 0 \quad \Rightarrow \quad \bar{u}_h \geq 0$$

$$\text{DMaxP} \quad \Leftrightarrow \quad \text{DMinP} \quad \Leftrightarrow \quad \text{DComP}$$

Definition (DMP)

$$\mathcal{T}_h \text{ fixed,} \quad f \geq 0 \ \& \ g \geq 0 \quad \Rightarrow \quad \bar{u}_h \geq 0 \quad (\text{everywhere in } \Omega)$$

Proof



$$\text{DMaxP} \Leftrightarrow \text{DMinP} \Leftrightarrow \text{DComP}$$

$$\text{DMaxP} \Leftrightarrow \text{DMinP:}$$

$$f \mapsto \bar{u}_h, \quad -f \mapsto -\bar{u}_h, \quad \min_{\bar{\Omega}} \bar{u}_h = - \max_{\bar{\Omega}} -\bar{u}_h$$



$$\text{DMaxP} \Leftrightarrow \text{DMinP} \Leftrightarrow \text{DComP}$$

$$\text{DMinP} \Rightarrow \text{DComP:}$$

$$f \geq 0, g \geq 0, \quad f \mapsto \bar{u}_h$$

$$\min_{\bar{\Omega}} \bar{u}_h \geq \min\{0, \min_{\partial\Omega} \bar{u}_h\} = \min\{0, \min_{\partial\Omega} g_h\} = 0$$

$$\text{DMaxP} \Leftrightarrow \text{DMinP} \Leftrightarrow \text{DComP}$$

$$\text{DComP} \Rightarrow \text{DMinP:}$$

$$f \geq 0, \quad f \mapsto \bar{u}_h, \quad g_h = u_h|_{\partial\Omega}, \quad \bar{g}_h = \min_{\partial\Omega} \bar{u}_h$$

$$\text{a) } \bar{g}_h \geq 0 \Rightarrow \bar{u}_h \geq 0$$

$$\min_{\Omega} \bar{u}_h \geq 0 = \min\{0, \min_{\partial\Omega} \bar{u}_h\}$$

$$\text{b) } \bar{g}_h < 0,$$

$$\tilde{f} = c\bar{g}_h \leq 0, \quad \tilde{g} = \bar{g}_h = \text{const.} \quad \mapsto \quad \tilde{u}_h = \bar{g}_h$$

$$f \geq 0 \geq \tilde{f}, \quad g_h \geq \tilde{g} \Rightarrow \bar{u}_h \geq \tilde{u}_h$$

$$\min_{\Omega} \bar{u}_h \geq \min_{\Omega} \tilde{u}_h = \bar{g}_h = \min_{\partial\Omega} \bar{u}_h = \min\{0, \min_{\partial\Omega} \bar{u}_h\}$$

\mathcal{T}_h ... N interior nodes, N^∂ boundary nodes, $\bar{N} = N + N^\partial$

$$\bar{u}_h(x) = \sum_{j=1}^{\bar{N}} \bar{y}_j \phi_j(x) = \sum_{j=1}^N y_j \phi_j(x) + \sum_{k=1}^{N^\partial} g_k \phi_{N+k}(x)$$

$$\bar{\mathbf{A}} \bar{\mathbf{y}} = \bar{\mathbf{F}} \quad \Leftrightarrow \quad \begin{bmatrix} \mathbf{A} & \mathbf{A}^\partial \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{g} \end{bmatrix} = \begin{bmatrix} \mathbf{F} \\ \mathbf{g} \end{bmatrix}$$

► $A_{ij} = a(\phi_j, \phi_i)$ $F_i = (f, \phi_i)$ $i = 1, 2, \dots, N, j = 1, 2, \dots, \bar{N}$.

► $\phi_j \geq 0$ $\sum_{j=1}^{\bar{N}} \phi_j \equiv 1$ $\bar{\mathbf{A}}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & -\mathbf{A}^{-1} \mathbf{A}^\partial \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$

Definition (DMP)

\mathcal{T}_h fixed, $f \geq 0$ & $g \geq 0 \Rightarrow \bar{u}_h \geq 0$ (everywhere in Ω)

Theorem (Ciarlet 1970)

DMP \Leftrightarrow (A1) $\bar{\mathbf{A}}^{-1} \geq 0$ (i.e. $\bar{\mathbf{A}}$ monotone)

(G1) $\xi + \mathbf{A}^{-1} \mathbf{A}^\partial \xi^\partial \geq 0$,

where $\xi = \underbrace{[1, \dots, 1]^\top}_{N \text{ times}}$ and $\xi^\partial = \underbrace{[1, \dots, 1]^\top}_{N^\partial \text{ times}}$

Remark ($g = 0$)

$f \geq 0 \Rightarrow \mathbf{F} \geq 0 \Leftrightarrow \underbrace{\mathbf{A}\mathbf{y} \geq 0 \Rightarrow \mathbf{y} \geq 0}_{\Leftrightarrow \mathbf{A} \text{ monotone}} \Leftrightarrow u_h \geq 0 \text{ in } \Omega$
 $\Leftrightarrow \mathbf{A}^{-1} \geq 0$

Monotone matrices



Let $\mathbf{A} \in \mathbb{R}^{N \times N}$.

- ▶ \mathbf{A} monotone if $\mathbf{A}\mathbf{y} \geq 0 \Rightarrow \mathbf{y} \geq 0$.
- ▶ \mathbf{A} M-matrix if $\text{off-diag}(\mathbf{A}) \leq 0$, \mathbf{A} nonsingular, and $\mathbf{A}^{-1} \geq 0$.
- ▶ \mathbf{A} Stieltjes if $\text{off-diag}(\mathbf{A}) \leq 0$ and \mathbf{A} s.p.d.

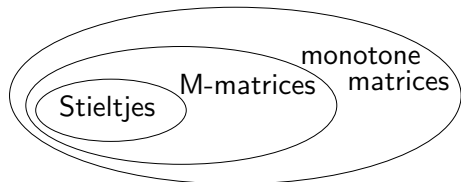
Theorem

Let \mathbf{A} be nonsingular.

\mathbf{A} monotone $\Leftrightarrow \mathbf{A}^{-1} \geq 0$.

Theorem (Varga 1962)

\mathbf{A} Stieltjes $\Rightarrow \mathbf{A}$ M-matrix.



Ph. Ciarlet conditions



- (a) $a_{ii} > 0, i = 1, \dots, N$
- (b) $a_{ij} \leq 0, i \neq j, i = 1, \dots, N, j = 1, \dots, N + N^\partial$
- (c) $\sum_{j=1}^{N+N^\partial} a_{ij} \geq 0, i = 1, \dots, N$
- (d) \mathbf{A} is irreducibly diagonally dominant

Theorem (Ciarlet 1970)

(a)–(d) \Rightarrow DMP

Proof.

(A1): Thm. (Varga): (a),(b),(d) $\Rightarrow \mathbf{A}^{-1} > 0$

$$(b) \Rightarrow -\mathbf{A}^{-1}\mathbf{A}^\partial = \mathbf{A}^{-1}(-\mathbf{A}^\partial) \geq 0 \Rightarrow \bar{\mathbf{A}}^{-1} \geq 0$$

(G1): (c) $\Leftrightarrow \mathbf{A}\xi + \mathbf{A}^\partial\xi^\partial \geq 0 \Rightarrow \xi + \mathbf{A}^{-1}\mathbf{A}^\partial\xi^\partial \geq 0 \quad \square$

$$(b) \quad a_{ij} \leq 0, \quad i \neq j, \quad i = 1, \dots, N, \quad j = 1, \dots, N + N^\partial$$

Theorem

$$(b) \Rightarrow DMP$$

Proof.

$$(A1): \mathbf{A} \text{ s.p.d. and } (b) \stackrel{\text{def}}{\Leftrightarrow} \text{Stieltjes matrix} \Rightarrow \text{M-matrix} \\ \Rightarrow \mathbf{A}^{-1} \geq 0$$

$$(b) \Rightarrow -\mathbf{A}^{-1}\mathbf{A}^\partial = \mathbf{A}^{-1}(-\mathbf{A}^\partial) \geq 0 \Rightarrow \bar{\mathbf{A}}^{-1} \geq 0$$

$$(G1): \sum_{j=1}^{N+N^\partial} a_{ij} = a\left(\sum_{j=1}^{N+N^\partial} \phi_j, \phi_i\right) = a(\mathbf{1}, \phi_i) = \int_{\Omega} c \phi_i \geq 0$$

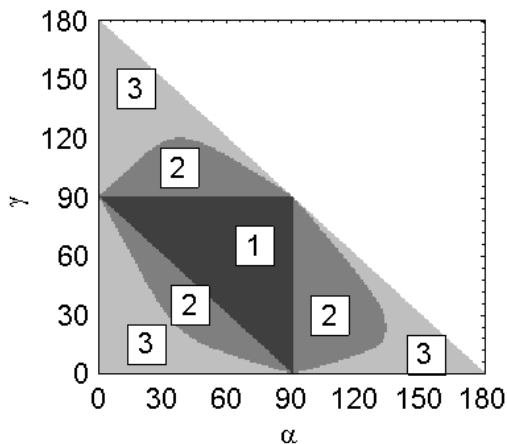
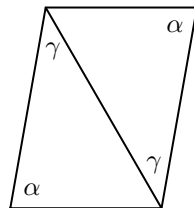
$$\Rightarrow (c) \Leftrightarrow \mathbf{A}\xi + \mathbf{A}^\partial \xi^\partial \geq 0 \Rightarrow \xi + \mathbf{A}^{-1}\mathbf{A}^\partial \xi^\partial \geq 0$$

□

Experiment I



$$-\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$



1 $\text{off-diag}(A) \leq 0$

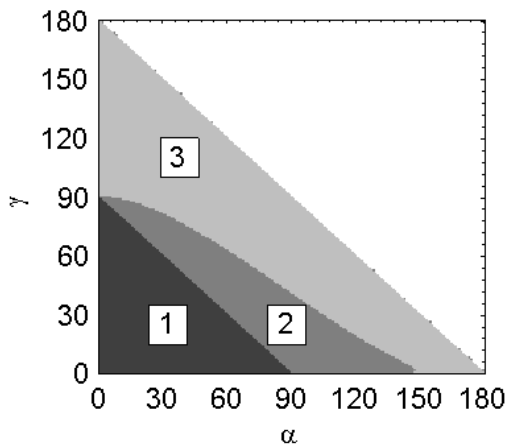
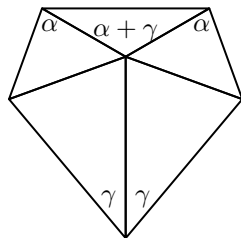
2 $\text{off-diag}(A) \not\leq 0, \mathbf{A}^{-1} \geq 0$

3 $\mathbf{A}^{-1} \not\geq 0$

Experiment II



$$-\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$



1 $\text{off-diag}(A) \leq 0$

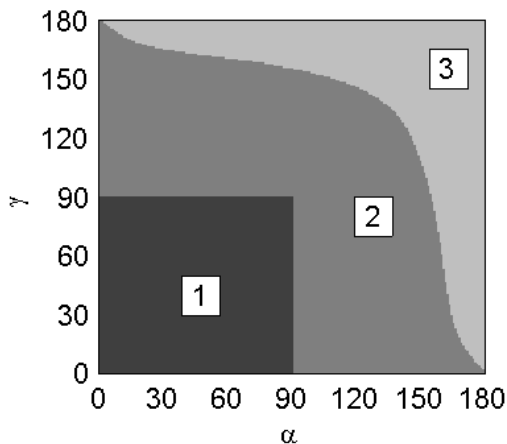
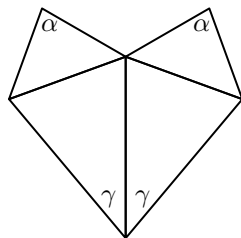
2 $\text{off-diag}(A) \not\leq 0, \mathbf{A}^{-1} \geq 0$

3 $\mathbf{A}^{-1} \not\geq 0$

Experiment III



$$-\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$



1 $\text{off-diag}(A) \leq 0$

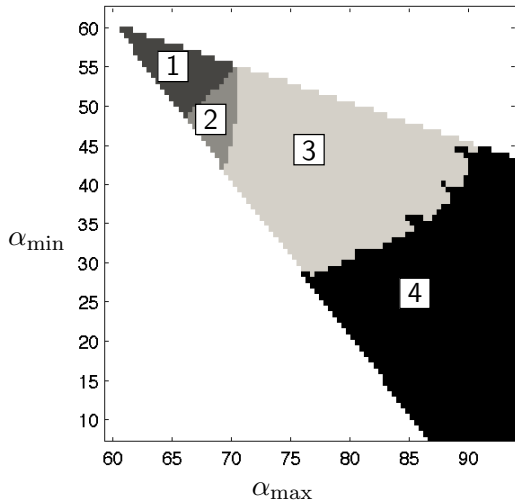
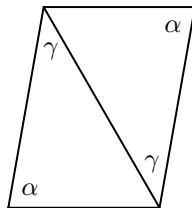
2 $\text{off-diag}(A) \not\leq 0, \mathbf{A}^{-1} \geq 0$

3 $\mathbf{A}^{-1} \not\geq 0$

Experiment prisms



$$-\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$



$$-\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$

$$\boxed{1} \quad \frac{1}{2} |T_{\max}| \tan \alpha_{\max}^{\mathcal{T}_h^G} \leq d_i^2 \leq |T_{\min}| \tan \alpha_{\min}^{\mathcal{T}_h^G} \Rightarrow \text{DMP}$$

$$\boxed{2} \quad d_L^{(P)} \leq d^{(P)} \leq d_U^{(P)} \quad \text{for all } P \in \mathcal{T}_{h,\tau} \Rightarrow \text{DMP}$$

$$d_L^{(P)} = \left(\frac{2 \cot \alpha_{\max}^{(T)}}{|T|} \right)^{-\frac{1}{2}}, \quad d_U^{(P)} = \left(\frac{\cot \alpha_{\text{med}}^{(T)} + \cot \alpha_{\min}^{(T)}}{2|T|} \right)^{-\frac{1}{2}}$$

$$\boxed{3} \quad \mathbf{A}^{-1} \geq 0 \Leftrightarrow \text{DMP}$$

$$\boxed{4} \quad \text{no DMP}$$

A. Hannukainen, S. Korotov, T. Vejchodský: Discrete maximum principle for FE solutions of the diffusion-reaction problem on prismatic meshes, accepted by J. Comput. Appl. Math., 2008.

Thank you for your attention

Tomáš Vejchodský¹ (vejchod@math.cas.cz)

Antti Hannukainen² (antti.hannukainen@hut.fi)

Sergey Korotov² (sergey.korotov@hut.fi)

¹ Institute of Mathematics, Academy of Sciences
Žitná 25, 115 67 Prague 1
Czech Republic

² Institute of Mathematics
Helsinki University of Technology
P.O. Box 1100, FIN-02015 Espoo
Finland

