



A vertex-centered and positivity-preserving scheme for anisotropic diffusion problems on arbitrary polygonal grids [☆]



Xiaoping Zhang^a, Shuai Su^b, Jiming Wu^{c,*}

^a School of Mathematics and Statistics and Computational Science Hubei Key Laboratory, Wuhan University, Wuhan, Hubei 430072, PR China

^b Graduate School of China Academy of Engineering Physics, Beijing, 100088, PR China

^c Institute of Applied Physics and Computational Mathematics, Beijing 100088, PR China

ARTICLE INFO

Article history:

Received 15 December 2016

Received in revised form 5 April 2017

Accepted 27 April 2017

Available online 4 May 2017

Keywords:

Diffusion equation

Vertex-centered scheme

Positivity-preserving

Nonlinear two-point flux approximation

ABSTRACT

We suggest a new positivity-preserving finite volume scheme for anisotropic diffusion problems on arbitrary polygonal grids. The scheme has vertex-centered, edge-midpoint and cell-centered unknowns. The vertex-centered unknowns are primary and have finite volume equations associated with them. The edge-midpoint and cell-centered unknowns are treated as auxiliary ones and are interpolated by the primary unknowns, which makes the final scheme a pure vertex-centered one. Unlike most existing positivity-preserving schemes, the construction of the scheme is based on a special nonlinear two-point flux approximation that has a fixed stencil and does not require the convex decomposition of the co-normal. In order to solve efficiently the nonlinear systems resulting from the nonlinear scheme, Picard method and its Anderson acceleration are discussed. Numerical experiments demonstrate the second-order accuracy and well positivity of the solution for heterogeneous and anisotropic problems on severely distorted grids. The high efficiency of the Anderson acceleration is also shown on reduction of the number of nonlinear iterations. Moreover, the proposed scheme does not have the so-called numerical heat-barrier issue suffered by most existing cell-centered and hybrid schemes. However, further improvements have to be made if the solution is very close to the machine precision and the mesh distortion is very severe.

© 2017 Elsevier Inc. All rights reserved.

1. Introduction

We consider the following diffusion problem

$$-\operatorname{div}(\Lambda \nabla u) = f \text{ in } \Omega, \quad (1)$$

$$u = g_D \text{ on } \Gamma_D, \quad (2)$$

$$-\Lambda \nabla u \cdot \mathbf{n} = g_N \text{ on } \Gamma_N, \quad (3)$$

where

[☆] This work was supported by the National Natural Science Foundation of China of China (Nos. 91330205, 11671313, 11571226).

* Corresponding author.

E-mail addresses: xpzhang.math@whu.edu.cn (X. Zhang), Sushuaiby@sina.cn (S. Su), wu_jiming@iapcm.ac.cn (J. Wu).

- (1) Ω is an open bounded connected polygonal domain in \mathbb{R}^2 ;
- (2) f is the source term, belonging to $L^2(\Omega)$;
- (3) Λ is a symmetric tensor such that (a) Λ is piecewise Lipschitz-continuous on Ω and (b) the set of eigenvalues of Λ is included in $[\lambda_{\min}, \lambda_{\max}]$ with $\lambda_{\min} > 0$;
- (4) $\partial\Omega = \bar{\Gamma}_D \cup \bar{\Gamma}_N$ is the boundary of Ω , \mathbf{n} denotes the exterior unit normal vector and g_D (resp. g_N) is the Dirichlet (resp. Neumann) boundary data defined on Γ_D (resp. Γ_N).

In addition, we assume that $f \geq 0$, $g_D \geq 0$ and $g_N \leq 0$, so that nonnegativity of the solution u is always assured.

Diffusion equations of the form (1) are at the core of complex models such as those in radiation hydrodynamics (RHD), magnetohydrodynamics (MHD), reservoir modeling, and so on. Finite volume (FV) discretizations in these applications are well received due to some important characteristics such as simplicity and local conservation. The readers are referred to, e.g., [6,8,9], for some recent developments. In many situations, one of the significant requirements of the FV schemes is that the discrete solution should be nonnegative. A scheme with such a property is usually referred to as positivity-preserving or monotone. In some applications such as RHD and reservoir simulations, the meshes are usually distorted and the media are heterogeneous and anisotropic, which may cause a scheme to violate the positivity-preserving property and exhibit nonphysical oscillations. Many scientists have made great efforts in the development of FV schemes that preserve the positivity-preserving property and at the same time, have approximately a second-order accuracy on severely distorted meshes in the case that the diffusion tensor is taken to be highly anisotropic, heterogeneous, and/or discontinuous.

As has been shown in [5,15,21] that, no linear consistent and conservative nine-point scheme can unconditionally respect the positivity-preserving property. One must therefore look for *nonlinear* schemes. In [22], C. Le Potier suggested a nonlinear two-point flux approximation to construct a cell-centered scheme on triangular meshes. Since then, this approach has been developed to obtain positivity-preserving schemes [3,14,16,26,35] or extreme-preserving ones [12,25] on general grids with general diffusion tensors. Two key ingredients of these works are the convex decomposition of the co-normal and the positivity-preserving interpolation of the auxiliary unknowns. By comparison, the later is quite difficult, even more difficult than the design of positivity-preserving FV scheme itself. The methods in [3,16,26,35] generally introduce auxiliary unknowns at mesh vertices or edge midpoints, whose values are interpolated from the neighboring cell-centered unknowns. When some of the weights in the interpolation formula are negative, the authors suggested that the corresponding auxiliary unknowns must be computed by some lower-order but positivity-preserving interpolation methods, such as the inverse distance weighting method [16]. Unfortunately, this technique usually leads to a great loss of accuracy on largely distorted meshes even in the case of a constant diffusion coefficient with a smooth solution. An interpolation-free approach based on local repositioning of cell centers was proposed in [17], but it can be used only for meshes with only one discontinuous interface per cell. A nonlinear two-point flux approximation a little different from that of Potier's was first proposed in [32] and further developed in [13] where both the positive interpolation of the auxiliary unknowns and the convex decomposition of the co-normal are not required. Therefore, it is possible for us to use arbitrary second-order interpolation algorithms and a relatively arbitrary co-normal decomposition. However, the truncation error in this approach is a little difficult to be analyzed rigorously.

To our knowledge, all the nonlinear positivity-preserving FV schemes are cell-centered except those in [23] and [6]. In [23], a nonlinear correction technique is applied to a class of hybrid schemes to insure the discrete maximum principle. Here we are more concerned about the nonlinear scheme in [6] where both the cell-centered and vertex-centered unknowns are treated as primary ones. This scheme is constructed on the primary mesh and its dual counterpart. Like most existing positivity-preserving schemes, it requires the convex decomposition of the co-normal. Moreover, it has FV equations for the two sets of primary unknowns so that no interpolation is involved. However, this nonlinear scheme can not reach the second-order accuracy in the case of discontinuous diffusion coefficients and the authors did not know how to improve it. We show in this work that the convex decomposition of the co-normal may result in the violation of the linearity-preserving property. To construct a vertex-centered nonlinear FV scheme with a second-order accuracy is among the main motivations of this paper.

Recently, the authors in [20] pointed out that many existing cell-centered and hybrid FV schemes for nonlinear parabolic equations, including the mimetic finite difference schemes [4,19], suffer the so-called heat-barrier issue. More explicitly, any schemes based on the harmonic averaging of cell-centered diffusion coefficients will break down when some of these coefficients go to zero or their ratio grows, which results in totally wrong numerical solution profiles in some strongly nonlinear problems such as the propagation of a nonlinear heat wave in a cold media. To our knowledge, the report of this issue can be traced back to [2]. Our numerical experiments indicate that all the cell-centered and hybrid linearity-preserving schemes we studied before, including those nonlinear ones in [12,13,32], suffer the same problem. The vertex-centered schemes suggested in [31,33] do not have the numerical heat-barrier issue but they are not positivity-preserving. To design a positivity-preserving scheme that can overcome the heat-barrier issue is another main motivation here.

The efficiency of nonlinear solvers is another important issue. At the present, in solving the nonlinear systems resulting from the aforementioned nonlinear FV schemes, Picard method or fixed-point iteration is frequently used because it preserves the M-matrix structure of the coefficient matrix, which is required to guarantee the positivity of the discrete solutions. In some extreme cases such as largely distorted meshes, Picard method can be very slow, which motivates us to seek an efficient nonlinear solver that does not spoil the M-matrix structure of the coefficient matrix. Anderson acceleration (AA) has been provided for nonlinear algebraic or transcendental equations in [1], and has been successfully used in

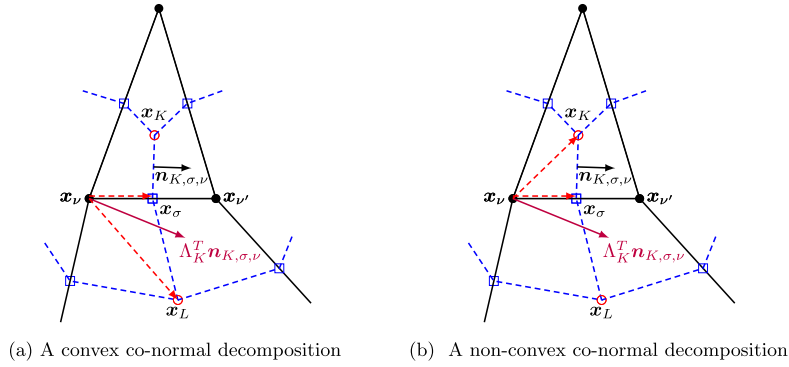


Fig. 2. Notations for the construction of one-sided flux.

their corresponding FV equations, instead, they are evaluated by primary unknowns, or Dirichlet boundary data if necessary, which will be discussed in details later. For simplicity of exposition, the formulation of the FV equation is illustrated for interior vertices. The discussion for vertices on Neumann boundary can be done analogously.

From now on, all the derivations are conducted under the following assumptions:

- The solution is smooth inside each primary cell and continuous on the whole domain Ω , while the diffusion tensor is piecewise constant with respect to the primary mesh;
- The possible discontinuities of the solution gradient and the diffusion tensor are only allowed to appear on the edges of the primary mesh.

Here we remark that the above two assumptions are the same as those in the derivation of cell-centered or hybrid FV schemes.

2.2. One-sided flux formulation

As usual, we define the flux vector as $\mathbf{F} = -\Lambda \nabla u$. Its normal component is then given by

$$\mathbf{F} \cdot \mathbf{n} = (-\Lambda \nabla u) \cdot \mathbf{n} = -\nabla u \cdot (\Lambda^T \mathbf{n}),$$

where \mathbf{n} is a unit vector normal to a certain primary or dual edge while $\Lambda^T \mathbf{n}$ is the so-called co-normal. The *one-sided flux* associated with each dual edge is constructed using only the information on one side of the dual edge. The one-sided fluxes are the bricks in building our vertex-centered and positivity-preserving scheme.

Let \mathbf{x}_v and $\mathbf{x}_{v'}$ be two generic interior vertices of the primary mesh shared by edge σ of the primary cell K , see Fig. 2. \mathbf{x}_σ , \mathbf{x}_K , \mathcal{E}_K and Λ_K denote the midpoint of σ , the center of K , the set of edges of K and the constant restriction of Λ on K , respectively. Denote by $\mathbf{n}_{K,\sigma,v}$ (resp. $\mathbf{n}_{K,\sigma,v'}$) the unit vector normal to dual edge $\mathbf{x}_\sigma \mathbf{x}_K$ outward to the dual cell associated with \mathbf{x}_v (resp. $\mathbf{x}_{v'}$). Obviously, $\mathbf{n}_{K,\sigma,v} = -\mathbf{n}_{K,\sigma,v'}$. In addition, let \mathcal{E}_v (resp. \mathcal{M}_v) be the set of primary edges (resp. primary cells) sharing \mathbf{x}_v .

2.2.1. An existing one-sided flux formulation

For the co-normal $\Lambda_K^T \mathbf{n}_{K,\sigma,v}$, since the dual cells are star-shaped, one can always find $L \in \mathcal{M}_v$ and $\tau \in \mathcal{E}_v \cap \mathcal{E}_L$, such that

$$|\mathbf{x}_\sigma \mathbf{x}_K| \Lambda_K^T \mathbf{n}_{K,\sigma,v} = \alpha_{L,\tau,v} (\mathbf{x}_L - \mathbf{x}_v) + \beta_{L,\tau,v} (\mathbf{x}_\tau - \mathbf{x}_v), \quad (4)$$

where

$$\alpha_{L,\tau,v} \geq 0, \quad \beta_{L,\tau,v} \geq 0, \quad \alpha_{L,\tau,v} + \beta_{L,\tau,v} > 0. \quad (5)$$

The geometric meaning of the above decomposition is that the co-normal is located between $\mathbf{x}_L - \mathbf{x}_v$ and $\mathbf{x}_\tau - \mathbf{x}_v$, see Fig. 2a where $\tau = \sigma$. (4) (with (5)) is the so-called *convex decomposition* of the co-normal, which is required in constructing most existing positivity-preserving schemes such as that in [6]. It follows from (4) that

$$\int_{\mathbf{x}_\sigma \mathbf{x}_K} \mathbf{F} \cdot \mathbf{n}_{K,\sigma,v} ds \approx \alpha_{L,\tau,v} (u(\mathbf{x}_v) - u(\mathbf{x}_L)) + \beta_{L,\tau,v} (u(\mathbf{x}_v) - u(\mathbf{x}_\tau)). \quad (6)$$

We claim that the one-sided flux approximation (6), derived from the convex decomposition (4), may not be linearity-preserving, i.e., its truncation error does not vanish if the solution is piecewise linear and the diffusion tensor is piecewise

constant with respect to the primary mesh. In order to see this clearly, let us take the case in Fig. 2a as an example where $\sigma \in \mathcal{E}_K \cap \mathcal{E}_L$ and $\tau = \sigma$. Assume further that

$$u = \begin{cases} a_K x + b_K y + c_K, & \text{on } K, \\ a_L x + b_L y + c_L, & \text{on } L. \end{cases}$$

Set $\mathbf{g}_K = (a_K, b_K)^T$ and $\mathbf{g}_L = (a_L, b_L)^T$. By direct calculations, we find that

$$\int_{\mathbf{x}_\sigma \mathbf{x}_K} \mathbf{F} \cdot \mathbf{n}_{K,\sigma,v} ds = [\alpha_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_L) + \beta_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_\sigma)] \cdot \mathbf{g}_K,$$

and

$$\alpha_{L,\sigma,v}(u(\mathbf{x}_v) - u(\mathbf{x}_L)) + \beta_{L,\sigma,v}(u(\mathbf{x}_v) - u(\mathbf{x}_\sigma)) = [\alpha_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_L) + \beta_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_\sigma)] \cdot \mathbf{g}_L.$$

The continuity of u across σ implies that

$$(\mathbf{x}_v - \mathbf{x}_\sigma) \cdot \mathbf{g}_K = (\mathbf{x}_v - \mathbf{x}_\sigma) \cdot \mathbf{g}_L.$$

Therefore, if (6) is linearity-preserving, we must have

$$\alpha_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_L) \cdot \mathbf{g}_K = \alpha_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_L) \cdot \mathbf{g}_L$$

or equivalently,

$$\alpha_{L,\sigma,v}(\mathbf{x}_v - \mathbf{x}_L) \cdot (\mathbf{g}_K - \mathbf{g}_L) = 0. \quad (7)$$

If the gradient of u happens to be discontinuous across σ so that $\mathbf{g}_K \neq \mathbf{g}_L$, then (7) will not hold in general, which explains why the scheme in [6] cannot achieve the second-order accuracy in the case of discontinuous diffusion coefficients.

2.2.2. A new one-sided flux formulation

In view of the discussion in the previous subsection, we suggest a new one-sided flux based on a non-convex decomposition of the co-normal. As shown in Fig. 2b, for the co-normal $\Lambda_K^T \mathbf{n}_{K,\sigma,v}$, we have the following fixed decomposition

$$|\mathbf{x}_\sigma \mathbf{x}_K| \Lambda_K^T \mathbf{n}_{K,\sigma,v} = \alpha_{K,\sigma,v}(\mathbf{x}_K - \mathbf{x}_v) + \beta_{K,\sigma,v}(\mathbf{x}_\sigma - \mathbf{x}_v), \quad (8)$$

where

$$\alpha_{K,\sigma,v} = \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)}{(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)}, \quad \beta_{K,\sigma,v} = \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(\mathbf{x}_K - \mathbf{x}_v)}{(\mathbf{x}_\sigma - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_K - \mathbf{x}_v)} \quad (9)$$

and \mathcal{R} denotes an operator that rotates a vector clockwise to its normal direction. Here we remark that one of the coefficients in (8) may be negative so that it is not a convex decomposition in general. It follows from (8) that

$$\int_{\mathbf{x}_\sigma \mathbf{x}_K} \mathbf{F} \cdot \mathbf{n}_{K,\sigma,v} ds \simeq \alpha_{K,\sigma,v}(u(\mathbf{x}_v) - u(\mathbf{x}_K)) + \beta_{K,\sigma,v}(u(\mathbf{x}_v) - u(\mathbf{x}_\sigma)),$$

here and hereafter \simeq indicates that the formula holds in the linearity-preserving sense. Then we obtain the following one-sided flux approximation with respect to dual edge $\mathbf{x}_\sigma \mathbf{x}_K$,

$$F_{K,\sigma,v} := \alpha_{K,\sigma,v}(u_v - u_K) + \beta_{K,\sigma,v}(u_v - u_\sigma), \quad (10)$$

where u_K, u_v and u_σ denote the approximations of u at $\mathbf{x}_K, \mathbf{x}_v$ and \mathbf{x}_σ , respectively.

Analogously, from the vector splitting

$$|\mathbf{x}_\sigma \mathbf{x}_K| \Lambda_K^T \mathbf{n}_{K,\sigma,v'} = \alpha_{K,\sigma,v'}(\mathbf{x}_K - \mathbf{x}_{v'}) + \beta_{K,\sigma,v'}(\mathbf{x}_\sigma - \mathbf{x}_{v'}), \quad (11)$$

we obtain another one-sided flux approximation,

$$F_{K,\sigma,v'} := \alpha_{K,\sigma,v'}(u_{v'} - u_K) + \beta_{K,\sigma,v'}(u_{v'} - u_\sigma), \quad (12)$$

where $u_{v'}$ denotes the approximation of u at $\mathbf{x}_{v'}$, and coefficients in (12) are computed in a way similar to those of (10).

Lemma 1. If the primary cell K is star-shaped with respect to its center \mathbf{x}_K , then for the coefficients in the one-sided fluxes (10) and (12), we have

$$\alpha_{K,\sigma,v} + \beta_{K,\sigma,v} = \alpha_{K,\sigma,v'} + \beta_{K,\sigma,v'} > 0, \quad (13)$$

$$\alpha_{K,\sigma,v} + \alpha_{K,\sigma,v'} = 0 \quad (14)$$

and

$$|\alpha_{K,\sigma,v}| \leq \frac{\lambda_{\max}}{\sin \angle \mathbf{x}_K \mathbf{x}_\sigma \mathbf{x}_v}, \quad (15)$$

where λ_{\max} denotes the maximal eigenvalue of Λ .

Proof. Since K is a star-shaped polygon with respect to \mathbf{x}_K , we have

$$(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v) = -(\mathbf{x}_\sigma - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_K - \mathbf{x}_v) = -|\mathbf{x}_\sigma \mathbf{x}_K| d_{K,\sigma,v}, \quad (16)$$

where $d_{K,\sigma,v}$ denotes the distance from the vertex \mathbf{x}_v to the dual edge $\mathbf{x}_\sigma \mathbf{x}_K$. Then, we deduce from (9) that

$$\alpha_{K,\sigma,v} + \beta_{K,\sigma,v} = \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_K)}{(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)} = \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathbf{n}_{K,\sigma,v}}{d_{K,\sigma,v}},$$

where we have used the identity $\mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_K) = -|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}$. Analogously, we have

$$\alpha_{K,\sigma,v'} + \beta_{K,\sigma,v'} = \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v'}^T \Lambda_K \mathbf{n}_{K,\sigma,v'}}{d_{K,\sigma,v'}},$$

where $d_{K,\sigma,v'}$ denotes the distance from the vertex $\mathbf{x}_{v'}$ to the dual edge $\mathbf{x}_\sigma \mathbf{x}_K$. Since \mathbf{x}_σ is the midpoint of σ , we have $d_{K,\sigma,v} = d_{K,\sigma,v'}$. Recalling that $\mathbf{n}_{K,\sigma,v'} = -\mathbf{n}_{K,\sigma,v}$ and Λ_K is positive definite, we immediately obtain (13). As for (14), we have

$$(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v) = -(\mathbf{x}_K - \mathbf{x}_{v'})^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_{v'})$$

and as a result,

$$\begin{aligned} \alpha_{K,\sigma,v} + \alpha_{K,\sigma,v'} &= \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)}{(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)} + \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v'}^T \Lambda_K \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_{v'})}{(\mathbf{x}_K - \mathbf{x}_{v'})^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_{v'})} \\ &= \frac{|\mathbf{x}_\sigma \mathbf{x}_K| \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(2\mathbf{x}_\sigma - \mathbf{x}_v - \mathbf{x}_{v'})}{(\mathbf{x}_K - \mathbf{x}_v)^T \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)} = 0. \end{aligned}$$

Finally, by (9) and (16), we have

$$\alpha_{K,\sigma,v} = -\frac{\mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)}{d_{K,\sigma,v}} = -\frac{1}{\sin \angle \mathbf{x}_K \mathbf{x}_\sigma \mathbf{x}_v} \mathbf{n}_{K,\sigma,v}^T \Lambda_K \mathbf{n}_{K,\sigma,v},$$

where $\mathbf{n}_{K,\sigma} = \mathcal{R}(\mathbf{x}_\sigma - \mathbf{x}_v)/|\mathbf{x}_\sigma - \mathbf{x}_v|$ is a unit vector. Recalling once again that Λ_K is positive definite, we obtain (15). \square

2.3. A nonlinear two-point flux approximation

In this subsection, we shall use the new one-sided fluxes defined previously to construct a unique flux approximation. We still use Fig. 2b for exposition. For the interior dual edge $\mathbf{x}_\sigma \mathbf{x}_K$, based on (10) and (12), we define

$$\tilde{F}_{K,\sigma,v} = \mu_{K,\sigma,v} F_{K,\sigma,v} - \mu_{K,\sigma,v'} F_{K,\sigma,v'}, \quad \tilde{F}_{K,\sigma,v'} = \mu_{K,\sigma,v'} F_{K,\sigma,v'} - \mu_{K,\sigma,v} F_{K,\sigma,v} \quad (17)$$

where $\tilde{F}_{K,\sigma,v}$ (resp. $\tilde{F}_{K,\sigma,v'}$) approximates the flux outward from the dual cell associated with \mathbf{x}_v (resp. $\mathbf{x}_{v'}$), and $\mu_{K,\sigma,v}$ and $\mu_{K,\sigma,v'}$ are two positive parameters, satisfying

$$\mu_{K,\sigma,v} + \mu_{K,\sigma,v'} = 1. \quad (18)$$

Obviously, across the dual edge $\mathbf{x}_\sigma \mathbf{x}_K$, we have the flux continuity or the local conservation condition

$$\tilde{F}_{K,\sigma,v} + \tilde{F}_{K,\sigma,v'} = 0. \quad (19)$$

Substituting (10) and (12) into the first equation of (17) and rearranging the terms, we have

$$\tilde{F}_{K,\sigma,v} = \mu_{K,\sigma,v} (\alpha_{K,\sigma,v} + \beta_{K,\sigma,v}) u_v - \mu_{K,\sigma,v'} (\alpha_{K,\sigma,v'} + \beta_{K,\sigma,v'}) u_{v'} + B_{K,\sigma}(u_h), \quad (20)$$

where

$$B_{K,\sigma}(u_h) = \mu_{K,\sigma,v'} \xi_{K,\sigma,v'}(u_h) - \mu_{K,\sigma,v} \xi_{K,\sigma,v}(u_h), \quad (21)$$

$$\xi_{K,\sigma,v}(u_h) = \alpha_{K,\sigma,v} u_K + \beta_{K,\sigma,v} u_\sigma, \quad \xi_{K,\sigma,v'}(u_h) = \alpha_{K,\sigma,v'} u_K + \beta_{K,\sigma,v'} u_\sigma \quad (22)$$

and u_h denotes the finite volume solution defined on Ω whose values at \mathbf{x}_K , \mathbf{x}_v and \mathbf{x}_σ are u_K , u_v and u_σ , respectively. In order to obtain a two-point flux approximation that leads to a positivity-preserving scheme, we have to choose $\mu_{K,\sigma,v}$ and manipulate $B_{K,\sigma}(u_h)$ properly. Firstly, we choose

$$\mu_{K,\sigma,v} = \begin{cases} 0.5, & \xi_{K,\sigma,v}(u_h) = \xi_{K,\sigma,v'}(u_h) = 0, \\ \frac{|\xi_{K,\sigma,v'}(u_h)|}{|\xi_{K,\sigma,v}(u_h)| + |\xi_{K,\sigma,v'}(u_h)|}, & \text{otherwise.} \end{cases} \quad (23)$$

Then, we split $B_{K,\sigma}(u_h)$ as

$$B_{K,\sigma}(u_h) = \frac{u_v}{u_v + \epsilon} B_{K,\sigma}^+(u_h) - \frac{u_{v'}}{u_{v'} + \epsilon'} B_{K,\sigma}^-(u_h) + B_{K,\sigma}^{\epsilon,\epsilon'}(u_h),$$

where

$$B_{K,\sigma}^+(u_h) = \frac{|B_{K,\sigma}(u_h)| + B_{K,\sigma}(u_h)}{2}, \quad B_{K,\sigma}^-(u_h) = \frac{|B_{K,\sigma}(u_h)| - B_{K,\sigma}(u_h)}{2}, \quad (24)$$

$$B_{K,\sigma}^{\epsilon,\epsilon'}(u_h) = \frac{\epsilon}{u_v + \epsilon} B_{K,\sigma}^+(u_h) - \frac{\epsilon'}{u_{v'} + \epsilon'} B_{K,\sigma}^-(u_h),$$

and ϵ and ϵ' are two small positive numbers. It follows from (20) that

$$\tilde{F}_{K,\sigma,v} = A_{K,\sigma,v}(u_h) u_v - A_{K,\sigma,v'}(u_h) u_{v'} + B_{K,\sigma}^{\epsilon,\epsilon'}(u_h), \quad (25)$$

where

$$A_{K,\sigma,v}(u_h) = \mu_{K,\sigma,v} (\alpha_{K,\sigma,v} + \beta_{K,\sigma,v}) + \frac{B_{K,\sigma}^+(u_h)}{u_v + \epsilon}, \quad (26)$$

$$A_{K,\sigma,v'}(u_h) = \mu_{K,\sigma,v'} (\alpha_{K,\sigma,v'} + \beta_{K,\sigma,v'}) + \frac{B_{K,\sigma}^-(u_h)}{u_{v'} + \epsilon'}.$$

Finally, by truncating $B_{K,\sigma}^{\epsilon,\epsilon'}(u_h)$ in (25) and by (17), we arrive at

$$\tilde{F}_{K,\sigma,v} = A_{K,\sigma,v}(u_h) u_v - A_{K,\sigma,v'}(u_h) u_{v'}, \quad \tilde{F}_{K,\sigma,v'} = A_{K,\sigma,v'}(u_h) u_{v'} - A_{K,\sigma,v}(u_h) u_v, \quad (27)$$

here, for simplicity, we still use $\tilde{F}_{K,\sigma,v}$ and $\tilde{F}_{K,\sigma,v'}$ to denote the unique flux approximations. Obviously, the local conservation condition (19) is still maintained. Moreover, if $u_h \geq 0$, then $A_{K,\sigma,v}(u_h)$ and $A_{K,\sigma,v'}(u_h)$ are apparently nonnegative by (13), which is the basis of the entire formulation to maintain positivity-preserving.

2.4. The final nonlinear scheme

Based on the new nonlinear flux approximations given by (27), a nonlinear FV equation, associated with a vertex \mathbf{x}_v belongs to $\Omega \setminus \Gamma_D$, can be constructed as follows

$$\sum_{K \in \mathcal{M}_v} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_v} \tilde{F}_{K,\sigma,v} = \int_{K_v^*} f \, dx dy, \quad (28)$$

where K_v^* denotes the dual cell associated with \mathbf{x}_v .

3. Analysis of the truncation error

Throughout this section, without special mention, all the quantities defined in the previous section and involved u_h must be understood as their counterparts by replacing u_h with u . To guarantee a second-order accuracy, a necessary condition is that the truncation term $B_{K,\sigma}^{\epsilon,\epsilon'}(u_h)$ should satisfy the condition below

$$|B_{K,\sigma}^{\epsilon,\epsilon'}(u)| < Ch^2, \quad \forall \sigma \in \mathcal{E}_K, \forall K \in \mathcal{M}, \quad (29)$$

where \mathcal{M} denotes the set of primary cells, $h = \max_{K \in \mathcal{M}} h_K$ and h_K is the diameter of cell K . In the following discussion, notation (u) will be dropped for simplicity of exposition whenever there is no confusion. We first provide the following lemma.

Lemma 2. Let $\xi_{K,\sigma,v}$ and $\xi_{K,\sigma,v'}$ be defined by (22) with u_h replaced by u . If $\xi_{K,\sigma,v}\xi_{K,\sigma,v'} < 0$, then we have

$$|B_{K,\sigma}| \leq |\alpha_{K,\sigma,v}| |u(\mathbf{x}_K) - u(\mathbf{x}_\sigma)| \quad (30)$$

and

$$B_{K,\sigma}^{\epsilon,\epsilon'} = \begin{cases} \frac{\epsilon'}{u(\mathbf{x}_{v'}) + \epsilon'} B_{K,\sigma}, & \xi_{K,\sigma,v} > 0 \text{ and } \xi_{K,\sigma,v'} < 0, \\ \frac{\epsilon}{u(\mathbf{x}_v) + \epsilon} B_{K,\sigma}, & \xi_{K,\sigma,v} < 0 \text{ and } \xi_{K,\sigma,v'} > 0. \end{cases} \quad (31)$$

Proof. From (22), we have

$$\begin{aligned} \xi_{K,\sigma,v} &= (\alpha_{K,\sigma,v} + \beta_{K,\sigma,v}) u(\mathbf{x}_\sigma) + \alpha_{K,\sigma,v} (u(\mathbf{x}_K) - u(\mathbf{x}_\sigma)), \\ \xi_{K,\sigma,v'} &= (\alpha_{K,\sigma,v'} + \beta_{K,\sigma,v'}) u(\mathbf{x}_\sigma) + \alpha_{K,\sigma,v'} (u(\mathbf{x}_K) - u(\mathbf{x}_\sigma)). \end{aligned}$$

Then, by Lemma 1, we obtain

$$\xi_{K,\sigma,v} - \xi_{K,\sigma,v'} = 2\alpha_{K,\sigma,v} (u(\mathbf{x}_K) - u(\mathbf{x}_\sigma)) \quad (32)$$

and

$$\xi_{K,\sigma,v}\xi_{K,\sigma,v'} = (\alpha_{K,\sigma,v} + \beta_{K,\sigma,v})^2 u^2(\mathbf{x}_\sigma) - \alpha_{K,\sigma,v}^2 (u(\mathbf{x}_K) - u(\mathbf{x}_\sigma))^2. \quad (33)$$

Since $\xi_{K,\sigma,v}\xi_{K,\sigma,v'} < 0$, by substituting (23) and (18) into (21), we get

$$B_{K,\sigma} = \frac{2\xi_{K,\sigma,v}\xi_{K,\sigma,v'}}{\xi_{K,\sigma,v} - \xi_{K,\sigma,v'}}.$$

Using (32) and (33), we have

$$|B_{K,\sigma}| = \frac{\alpha_{K,\sigma,v}^2 (u(\mathbf{x}_K) - u(\mathbf{x}_\sigma))^2 - (\alpha_{K,\sigma,v} + \beta_{K,\sigma,v})^2 u^2(\mathbf{x}_\sigma)}{|\alpha_{K,\sigma,v}| |u(\mathbf{x}_K) - u(\mathbf{x}_\sigma)|}$$

which verifies (30). (31) is a natural consequence of the sign symbol of $B_{K,\sigma}$ and the definition of $B_{K,\sigma}^{\epsilon,\epsilon'}$. \square

Before the presentation of the main result, we need the following assumptions.

- (H1) $u \geq 0$ in Ω and $u \in C^2(\bar{K})$, $\forall K \in \mathcal{M}$;
- (H2) All zero points of u , if exist, are extreme points;
- (H3) There exists a positive constant γ , independent of h , such that

$$\sin \angle \mathbf{x}_K \mathbf{x}_\sigma \mathbf{x}_v \geq \gamma, \quad \forall \sigma \in \mathcal{E}_K, \quad \forall K \in \mathcal{M}.$$

- (H4) When $u(\mathbf{x}_v)u(\mathbf{x}_{v'}) \neq 0$, the positive numbers ϵ and ϵ' are chosen in such a way that

$$\epsilon \leq \frac{h}{1-h} u(\mathbf{x}_v), \quad \epsilon' \leq \frac{h}{1-h} u(\mathbf{x}_{v'}),$$

where $0 < h < 1$.

Theorem 1. Assume that $\mu_{K,\sigma,v}$ and $\mu_{K,\sigma,v'}$ are specified by (23) and (18). Then, under the assumptions (H1), (H2), (H3) and (H4), $B_{K,\sigma}^{\epsilon,\epsilon'}$ satisfies (29).

Proof. We first consider the first case where $\xi_{K,\sigma,v}\xi_{K,\sigma,v'} \geq 0$. In this case, we always have $B_{K,\sigma} = 0$ so that $B_{K,\sigma}^{\epsilon,\epsilon'} = 0$. What remains is to prove the case where $\xi_{K,\sigma,v}\xi_{K,\sigma,v'} < 0$. By Lemma 2 and (H1), we have

$$|B_{K,\sigma}^{\epsilon,\epsilon'}| \leq |\alpha_{K,\sigma,v}| |u(\mathbf{x}_\sigma) - u(\mathbf{x}_K)| \max \left(\frac{\epsilon}{u(\mathbf{x}_v) + \epsilon}, \frac{\epsilon'}{u(\mathbf{x}_{v'}) + \epsilon'} \right). \quad (34)$$

From (H3) and Lemma 1, $|\alpha_{K,\sigma,v}|$ has an upper bound independent of h . Therefore, if $u(\mathbf{x}_v)u(\mathbf{x}_{v'}) \neq 0$, one can deduce from (34) and (H4) that

$$|B_{K,\sigma}^{\epsilon,\epsilon'}| \leq |\alpha_{K,\sigma,v}| |u(\mathbf{x}_\sigma) - u(\mathbf{x}_K)| h \leq Ch^2.$$

If $u(\mathbf{x}_v)u(\mathbf{x}_{v'}) = 0$, without loss of generality, we only consider that \mathbf{x}_v is the zero point of u , i.e., $u(\mathbf{x}_v) = 0$. Then, from (H2), we find that $\nabla u(\mathbf{x}_v) = 0$. Note that $u \in C^2(\bar{K})$. By Taylor expansion, we have

$$|u(\mathbf{x}_\sigma) - u(\mathbf{x}_K)| = |\nabla u(\xi) \cdot (\mathbf{x}_\sigma - \mathbf{x}_K)| \leq \|\nabla u(\xi) - \nabla u(\mathbf{x}_v)\| \|\mathbf{x}_\sigma - \mathbf{x}_K\| \leq Ch^2,$$

where $\xi = \mathbf{x}_K + \theta(\mathbf{x}_\sigma - \mathbf{x}_K)$ with $\theta \in (0, 1)$. Substituting this into (34), we reach

$$|B_{K,\sigma}^{\epsilon,\epsilon'}| \leq |\alpha_{K,\sigma,v}| |u(\mathbf{x}_\sigma) - u(\mathbf{x}_K)| \leq Ch^2,$$

which completes the proof. \square

Remark 1. When $u(\mathbf{x}_v)$ or $u(\mathbf{x}_{v'})$ is very close to the machine precision, (H4) may be violated and the truncation error reduces to $O(h)$. In order to decrease the possibility of this violation, ϵ and ϵ' are chosen to be very small numbers in practical computation, say, 10^{-10} in double precision.

4. The interpolation of the auxiliary unknowns

In design of both vertex-centered and cell-centered FV schemes, the auxiliary unknowns must be interpolated by primary unknowns. A desirable interpolation algorithm is usually required to be second-order, simple, positivity-preserving, topology-independent, discontinuity-independent and so on, which is quite difficult for cell-centered schemes. As for the present scheme, since the auxiliary unknowns are defined at the centers and edge midpoints of the primary cells, they can be easily interpolated by the vertex unknowns. Let \mathbf{x}_σ be the midpoint of σ whose endpoints are \mathbf{x}_v and $\mathbf{x}_{v'}$, \mathbf{x}_K is the geometric center of K . Then, for auxiliary unknowns defined at \mathbf{x}_σ and \mathbf{x}_K , we have

$$u_\sigma = \frac{1}{2}(u_v + u_{v'}), \quad u_K = \frac{1}{n_K} \sum_{\mathbf{x}_v \in \mathcal{V}_K} u_v, \quad (35)$$

where \mathcal{V}_K (resp. n_K) denotes the set (resp. number) of vertices in K . Obviously, this interpolation algorithm satisfies all the aforementioned good properties, which is a major advantage of the present vertex-centered scheme over the interpolation based cell-centered schemes.

In the case where \mathbf{x}_K is not the geometric center of K , the second interpolation formula in (35) will generally not have the second-order accuracy so that we have to seek a different one. From (26), we can see that the present vertex-centered scheme does not require a positivity-preserving interpolation for the auxiliary unknowns. Hence, we have many choices to design a second-order interpolation algorithm that is not necessarily positivity-preserving, for example, the least square interpolation is a desirable candidate. However, in the case of nonlinear diffusion coefficients, such as $\Lambda(u) \propto u^{5/2}$ or $\Lambda(u) \propto u^3$, the positivity-preserving interpolation of the cell-centered unknowns is still a requirement. In this case, we adopt the second-order and positivity-preserving interpolation algorithm suggested in [31]. Here, for self-completeness of the present paper, we give a brief description. Let the vertices of cell K be denoted as \mathbf{x}_i ($1 \leq i \leq n_K$) which are ordered anticlockwisely. We have the formula below

$$u_K = \frac{1}{\sum_{i=1}^{n_K} (\zeta_i + \eta_i)} \sum_{i=1}^{n_K} (\zeta_i u_{j(i)} + \eta_i u_{j(i)+1}), \quad (36)$$

where u_l ($l = j(i), j(i) + 1$) denotes the approximation of u at \mathbf{x}_l ,

$$\begin{aligned} \zeta_i &= \frac{\mathcal{R}(\mathbf{x}_{i+1} - \mathbf{x}_i)^T \mathcal{R}(\mathbf{x}_{j(i)+1} - \mathbf{x}_K)}{(\mathbf{x}_{j(i)} - \mathbf{x}_K)^T \mathcal{R}(\mathbf{x}_{j(i)+1} - \mathbf{x}_K)}, \\ \eta_i &= \frac{\mathcal{R}(\mathbf{x}_{i+1} - \mathbf{x}_i)^T \mathcal{R}(\mathbf{x}_{j(i)} - \mathbf{x}_K)}{(\mathbf{x}_{j(i)+1} - \mathbf{x}_K)^T \mathcal{R}(\mathbf{x}_{j(i)} - \mathbf{x}_K)} \end{aligned} \quad (37)$$

and index $j(i)$ is selected in such a way that the normal vector $\mathcal{R}(\mathbf{x}_{i+1} - \mathbf{x}_i)$ is located between $\mathbf{x}_{j(i)} - \mathbf{x}_K$ and $\mathbf{x}_{j(i)+1} - \mathbf{x}_K$.

5. Monotonicity and iteration of the nonlinear scheme

The nonlinear scheme (28) can be formulated as the matrix form below

$$\mathbb{M}(\mathbf{U})\mathbf{U} = \mathbf{F}(\mathbf{U}), \quad (38)$$

where \mathbf{U} denotes the unknown vector and $\mathbb{M}(\mathbf{U})$ the coefficient matrix. The right-hand side vector $\mathbf{F}(\mathbf{U})$ is generated by the source term and the boundary data.

Algorithm 1: Picard method.

1 Choose a small positive value ϵ_{non} and $\mathbf{U}^0 \geq 0$.
 2 **for** $k = 1, 2, \dots$ **do**
 3 Solve the linear system

$$\mathbb{M}(\mathbf{U}^k)\mathbf{U}^{k+1} = \mathbf{F}(\mathbf{U}^k). \quad (39)$$

4 Use

$$\|\mathbb{M}(\mathbf{U}^{k+1})\mathbf{U}^{k+1} - \mathbf{F}(\mathbf{U}^{k+1})\| \leq \epsilon_{non} \|\mathbb{M}(\mathbf{U}^0)\mathbf{U}^0 - \mathbf{F}(\mathbf{U}^0)\| \quad (40)$$

as the convergence criterion.

5 **end**

5.1. The fixed-point iteration

The nonlinear system can be solved by the fixed-point or Picard iteration.

Theorem 2. Let $f \geq 0$, $g_D \geq 0$, and $g_N \leq 0$. If $\mathbf{U}^0 \geq 0$ and linear systems in the Picard iterations are solved exactly, then $\mathbf{U}^k \geq 0$ for $k \geq 1$.

Proof. There are similar conclusions in a number of papers, for instance, [6,13,17,35]. Here we just provide a sketch. Actually, we should prove that $\mathbf{U}^{k+1} \geq 0$ if $\mathbf{U}^k \geq 0$, which in turn reduces to prove $\mathbb{M}^{-1}(\mathbf{U}^k) \geq 0$ and $\mathbf{F}(\mathbf{U}^k) \geq 0$. In the following we shall drop (\mathbf{U}^k) or $(u_h^{(k)})$ if there is no confusion. By (13) and noticing $\mathbf{U}^k \geq 0$, we find that the coefficients $A_{K,\sigma,v}$ and $A_{K,\sigma,v'}$ in the nonlinear two-point flux approximation (27) are positive. Hence it follows from (27) and (28) that matrix \mathbb{M} has the properties listed below:

- \mathbb{M} is irreducible;
- All diagonal entries of matrix \mathbb{M} are positive;
- All off-diagonal entries of \mathbb{M} are non-positive;
- Each column sum in \mathbb{M} is non-negative and there exists a column with a positive sum.

Then we deduce that \mathbb{M}^T is an M-matrix. As a consequence, \mathbb{M}^{-T} and in turn, \mathbb{M}^{-1} are nonnegative. Under the assumption of $f \geq 0$, $g_D \geq 0$, and $g_N \leq 0$, one can deduce that $\mathbf{F} \geq 0$. This ends the proof. \square

5.2. The Anderson acceleration (AA)

Let m be a fixed positive integer. The Anderson acceleration of the Picard method is given in the following algorithm, see, e.g., [18] for details.

In practice, the solution of the constrained minimization problem (41) is shifted to the solution of a saddle-point problem [18]. The choice of m is another important issue, which will affect the efficiency of the iteration significantly. It is likely a problem-dependent issue and in the experiments reported in Section 6, the effective choice of m ranges from 2 to 11.

6. Numerical experiments

In this section, we examine the numerical performance of the vertex-centered and positive-preserving scheme (VPPS) discussed in the previous sections. The scheme with the Picard (resp. Anderson acceleration) method is denoted as VPPS-P (resp. VPPS-AA) for simplicity. We investigate the discrete relative errors for both the solution and the flux:

$$E_u = \left(\sum_{\mathbf{x}_v \in \bar{\Omega}} |K_v^*| (u_v - u(\mathbf{x}_v))^2 / \sum_{\mathbf{x}_v \in \bar{\Omega}} |K_v^*| u^2(\mathbf{x}_v) \right)^{1/2},$$

$$E_q = \left(\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} S_\sigma \left(\tilde{F}_{K,\sigma,v} - \tilde{F}_{K,\sigma,v}^{ext} \right)^2 / \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} S_\sigma \left(\tilde{F}_{K,\sigma,v}^{ext} \right)^2 \right)^{1/2},$$

where S_σ is an area associated with σ (for example, $S_\sigma = |K|/n_K$), and $\tilde{F}_{K,\sigma,v}^{ext} (= \tilde{F}_{K,\sigma,v'}^{ext})$ denotes the exact flux across the dual edge $\mathbf{x}_K \mathbf{x}_\sigma$. The rate of convergence R_α ($\alpha = u, q$) is obtained by a least squares fit on the ones computed on each two successive meshes by the following formula

$$R_\alpha = \frac{\log[E_\alpha(h_2)/E_\alpha(h_1)]}{\log(h_2/h_1)},$$

Algorithm 2: Anderson acceleration of Picard method.

- 1 Choose a small positive value ϵ_{non} and $\mathbf{U}^0 \geq 0$.
- 2 Apply two Picard iterations (39) and set $\tilde{\mathbf{U}}^k = \mathbf{U}^k$, $\delta \mathbf{U}^k = \tilde{\mathbf{U}}^k - \mathbf{U}^{k-1}$, $k = 1, 2$.
- 3 **for** $k = 2, \dots$ **do**
- 4 $m_k = \min(m, k)$
- 5 Determine weights $\alpha_1, \dots, \alpha_{m_k}$ by solving the minimization problem

$$\min \left\| \sum_{i=1}^{m_k} \alpha_i \delta \mathbf{U}^{k-m_k+i} \right\| \quad (41)$$

subjected to the constraint

$$\sum_{i=1}^{m_k} \alpha_i = 1. \quad (42)$$

- 6 Set new iterate

$$\mathbf{U}^{k+1} = \sum_{i=1}^{m_k} \alpha_i \tilde{\mathbf{U}}^{k-m_k+i}. \quad (43)$$

- 7 Lift $\mathbf{U}^{k+1} := \mathbf{U}^{k+1} - \min(\mathbf{U}^{k+1}, 0) \mathbf{e}$ where \mathbf{e} is a vector with all entries equal to 1.
- 8 Use formula (40) as the convergence criterion.
- 9 Solve the linear system

$$\mathbb{M}(\mathbf{U}^{k+1}) \tilde{\mathbf{U}}^{k+1} = \mathbf{F}(\mathbf{U}^{k+1})$$

and set $\delta \mathbf{U}^{k+1} = \tilde{\mathbf{U}}^{k+1} - \mathbf{U}^{k+1}$.

10 **end**

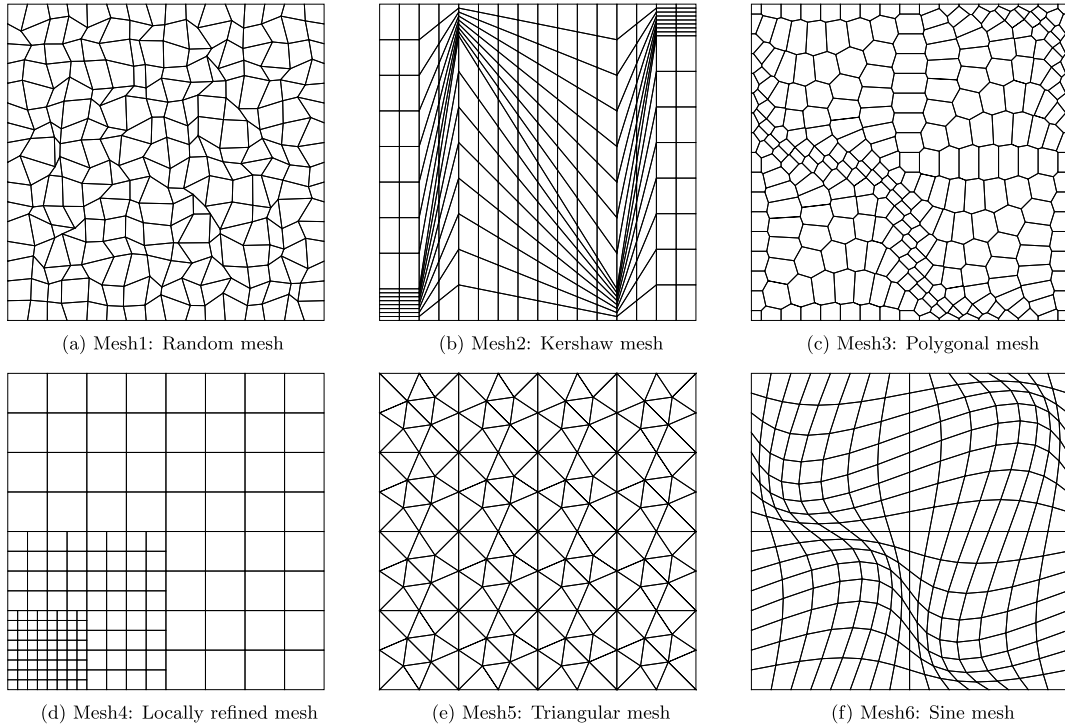


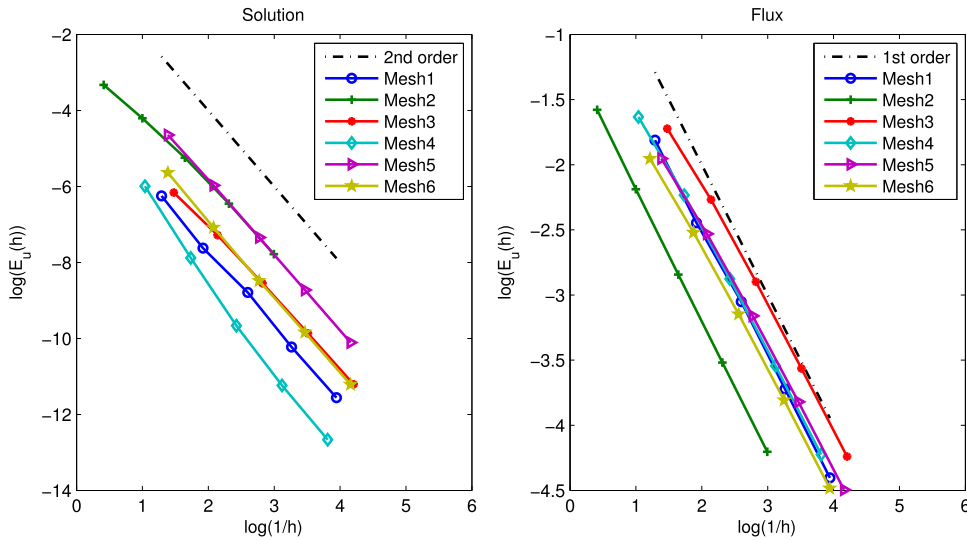
Fig. 3. Six mesh types used in the numerical tests.

where h_1, h_2 denote the mesh sizes of the two successive meshes, and $E_\alpha(h_1), E_\alpha(h_2)$ denote the corresponding discrete errors. In addition, we choose $\epsilon_{non} = 10^{-8}$ and $\epsilon = \epsilon' = 10^{-10}$ throughout this section.

Six mesh types are presented in Fig. 3, which are used in the following tests. The values of γ defined in (H3) on each level of the six mesh types are presented in Table 1. One can see that these six mesh types satisfy (H3).

Table 1The values of γ in (H3) among the six mesh types.

Level	Mesh1	Mesh2	Mesh3	Mesh4	Mesh5	Mesh6
1	0.502	0.204	0.772	0.848	0.928	0.326
2	0.281	0.204	0.690	0.848	0.928	0.312
3	0.137	0.204	0.382	0.848	0.928	0.311
4	0.125	0.204	0.198	0.848	0.928	0.310
5	0.046	0.204	0.036	0.848	0.928	0.310

**Fig. 4.** Convergence results of VPPS.

In addition, in order to verify (H4), we denote

$$\varrho := \min_{u_v \neq 0} \frac{h}{1-h} u_v \quad (44)$$

and evaluate it on each mesh level in the last step of the nonlinear iteration.

6.1. Smooth diffusion tensor

We first consider the homogeneous, mildly anisotropic problem investigated in [11]. The diffusion tensor is

$$\Lambda = \begin{pmatrix} 1.5 & 0.5 \\ 0.5 & 1.5 \end{pmatrix}$$

and the exact solution is given by

$$u(x, y) = \frac{1}{2} \left[\frac{\sin((1-x)(1-y))}{\sin(1)} + (1-x)^3(1-y)^2 \right],$$

while the boundary conditions and source term are determined accordingly. We solve this problem on the unit square domain $[0, 1]^2$ and use a sequence of six mesh types, see Fig. 3, in this numerical test.

The numerical results of VPPS are given respectively in Table 2 and graphically depicted in Fig. 4 as log-log plots of the discrete errors versus the characteristic mesh size h . The optimal convergence results are observed on all meshes. From Table 2, we can see that $\varrho > \varepsilon(\varepsilon')$ is always maintained on all mesh levels, which means that (H4) is valid.

In addition, we also investigate the performance of AA method. Table 3 presents the minimum solution by Picard method and the number of nonlinear iterations in the Picard and AA methods with various values of m on the last mesh refinement level. These results demonstrate that if the number of nonlinear iterations in Picard method is large, the AA method reduces it significantly (up to 3 times). These results suggest that $m = 7$ is a reasonable choice for this type of problem since the number of nonlinear iterations is small and the overhead caused by the implementation of the AA method is negligible.

At last, we give a test on Shestakov mesh, see Fig. 5. We find that VPPS fails at the fourth mesh level where $|B_{K,\sigma}^{\varepsilon,\varepsilon'}|/h^2$ increases largely in the Picard iteration process.

Table 2

The numerical results of VPPS on six mesh types.

Mesh level		1	2	3	4	5
Mesh1	E_u	1.931E-03	4.882E-04	1.520E-04	3.624E-05	9.555E-06
	Ratio		2.181	1.718	2.158	1.957
	E_q	1.634E-01	8.670E-02	4.733E-02	2.421E-02	1.223E-02
	Ratio		1.005	0.891	1.009	1.002
	ϱ	3.461E-03	2.642E-04	3.466E-05	7.223E-06	8.195E-07
Mesh2	E_u	3.586E-02	1.492E-02	5.280E-03	1.575E-03	4.174E-04
	Ratio		1.493	1.613	1.807	1.948
	E_q	2.064E-01	1.121E-01	5.823E-02	2.964E-02	1.494E-02
	Ratio		1.040	1.016	1.009	1.005
	ϱ	7.741E-10	3.069E-04	4.550E-05	5.950E-06	7.472E-07
Mesh3	E_u	2.110E-03	6.888E-04	1.943E-04	5.186E-05	1.360E-05
	Ratio		1.687	1.852	1.913	1.933
	E_q	1.786E-01	1.036E-01	5.510E-02	2.833E-02	1.442E-02
	Ratio		0.821	0.924	0.963	0.975
	ϱ	2.321E-05	3.607E-05	7.525E-06	6.456E-07	1.511E-07
Mesh4	E_u	2.495E-03	3.781E-04	6.339E-05	1.318E-05	3.155E-06
	Ratio		2.722	2.576	2.265	2.063
	E_q	1.952E-01	1.072E-01	5.632E-02	2.889E-02	1.464E-02
	Ratio		0.865	0.928	0.963	0.981
	ϱ	2.057E-02	1.996E-03	2.251E-04	2.683E-05	3.278E-06
Mesh5	E_u	9.590E-03	2.547E-03	6.471E-04	1.626E-04	4.071E-05
	Ratio		1.913	1.976	1.993	1.998
	E_q	1.419E-01	7.956E-02	4.241E-02	2.193E-02	1.116E-02
	Ratio		0.835	0.908	0.951	0.975
	ϱ	4.625E-03	4.883E-04	5.666E-05	6.841E-06	8.410E-07
Mesh6	E_u	3.553E-03	8.419E-04	2.070E-04	5.341E-05	1.350E-05
	Ratio		2.181	2.048	1.961	1.986
	E_q	1.417E-01	8.039E-02	4.295E-02	2.220E-02	1.129E-02
	Ratio		0.858	0.915	0.955	0.977
	ϱ	9.871E-04	2.109E-04	3.484E-05	4.979E-06	6.629E-07

Table 3

The minimum solution and the number of nonlinear iterations.

Mesh	VPPS-P		VPPS-AA					
	U_{\min}	$m = 1$	$m = 2$	$m = 3$	$m = 5$	$m = 7$	$m = 9$	$m = 11$
Mesh1	4.24E-05	21	19	16	15	15	15	15
Mesh2	1.42E-05	103	78	66	43	41	38	37
Mesh3	9.98E-06	17	15	15	15	15	15	15
Mesh4	1.45E-04	12	11	11	12	12	13	13
Mesh5	5.30E-05	10	9	9	9	9	9	9
Mesh6	3.35E-05	110	53	39	27	26	26	25

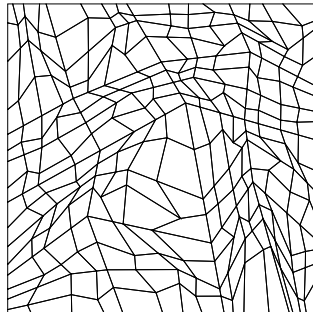
**Fig. 5.** Shestakov mesh.

Table 4

The numerical results of VPPS on Mesh1.

Mesh level	VPPS				
	E_u	Ratio	E_q	Ratio	ϱ
1	1.042E–03		1.049E+00		8.005E–01
2	2.968E–04	1.873	5.781E–01	0.888	2.464E–01
3	9.499E–05	1.783	2.991E–01	1.031	9.767E–02
4	2.788E–05	1.845	1.533E–01	1.006	4.492E–02
5	7.692E–06	1.890	7.724E–02	1.006	2.100E–02

Table 5

The number of nonlinear iterations for VPPS on Mesh1.

Mesh level	VPPS-P	VPPS-AA					
	$m = 1$	$m = 2$	$m = 3$	$m = 5$	$m = 7$	$m = 9$	$m = 11$
1	29	19	18	15	14	14	14
2	47	33	25	24	21	22	22
3	69	51	38	32	29	29	29
4	96	51	39	37	35	35	35
5	130	94	66	45	44	41	41

6.2. Discontinuous diffusion tensor

We deal with the problem (1)–(2) on $\Omega = [0, 1]^2$, and choose

$$\Lambda = \begin{cases} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, & x \leq 0.5, \\ \begin{pmatrix} 10 & 3 \\ 3 & 1 \end{pmatrix}, & x > 0.5. \end{cases}$$

The exact solution is

$$u(x, y) = \begin{cases} 1 - 2y^2 + 4xy + 6x + 2y, & x \leq 0.5, \\ -2y^2 + 1.6xy - 0.6x + 3.2y + 4.3, & x > 0.5. \end{cases}$$

We employ Mesh1 and all the mesh nodes located on the line $x = 0.5$ were distorted only in the y -direction.

The discrete relative errors of solution, flux and ϱ for scheme VPPS on each refinement level are given in Table 4, and the numbers of nonlinear iterations of VPPS-P and VPPS-AA with various values of m are presented in Table 5. From Table 4, one can see that the solution error is of second-order. By comparison, the NLMDDFV method [6] in this example is only first-order accurate. In addition, ϱ is far greater than ϵ and ϵ' we choose, which means the assumption (H4) is maintained. From Table 5, we can see that the AA method exhibits smaller dependence on the mesh size compared to the Picard method.

6.3. Heterogeneous rotating anisotropy

Problem (1)–(2) is now defined in $\Omega = [0, 1]^2$ with a rotating anisotropic diffusion tensor:

$$\Lambda = \frac{1}{x^2 + y^2} \begin{pmatrix} \beta x^2 + y^2 & (\beta - 1)xy \\ (\beta - 1)xy & x^2 + \beta y^2 \end{pmatrix},$$

where β characterizes the level of anisotropy. We consider the smooth exact solution $u(x, y) = \sin(\pi x) \sin(\pi y)$ and use the uniform square mesh and the unstructured triangular mesh Mesh5 with 5 mesh levels in this test. Obviously, we can find that $\gamma \equiv 1.0$ in each level on the uniform square mesh.

For various anisotropy $\beta = 1, 10^{-3}$ and 10^{-6} , the convergence rate on both two families of meshes are depicted in Fig. 6. We observe that VPPS achieves the optimal rate of convergence, second-order for the solution and first-order for the flux. In Table 6, we present the minimum solution and the number of nonlinear iterations, which shows the efficiency of the AA method.

6.4. Positivity of the discrete solution

This test problem, inspired by [16,18], is a highly anisotropic one, which aims to test the positivity-preserving property of FV schemes. The computational domain is a unit square with a square hole in the center, i.e., $\Omega = [0, 1]^2 \setminus [4/9, 5/9]^2$. The

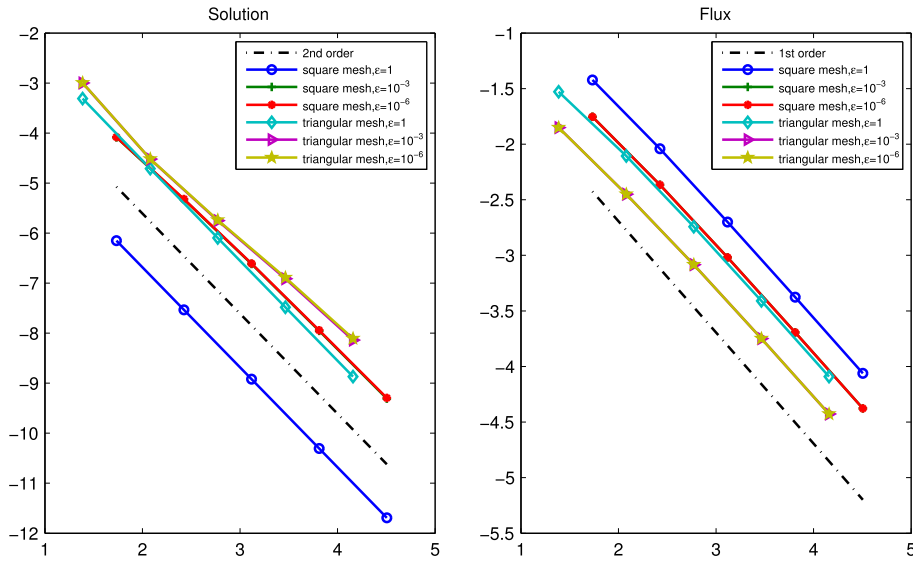


Fig. 6. Convergence results of VPPS.

Table 6

The minimal solution and the number of nonlinear iterations on Mesh5 with $\beta = 10^{-3}$.

Mesh level	VPPS-P		VPPS-AA					
	U_{\min}	$m = 1$	$m = 2$	$m = 3$	$m = 5$	$m = 7$	$m = 9$	$m = 11$
1	1.80E-01	47	20	18	14	13	13	13
2	4.70E-02	61	28	26	22	22	22	22
3	1.19E-02	121	50	49	31	31	32	30
4	2.98E-03	189	74	69	39	38	38	38
5	7.45E-04	257	77	80	46	45	45	43

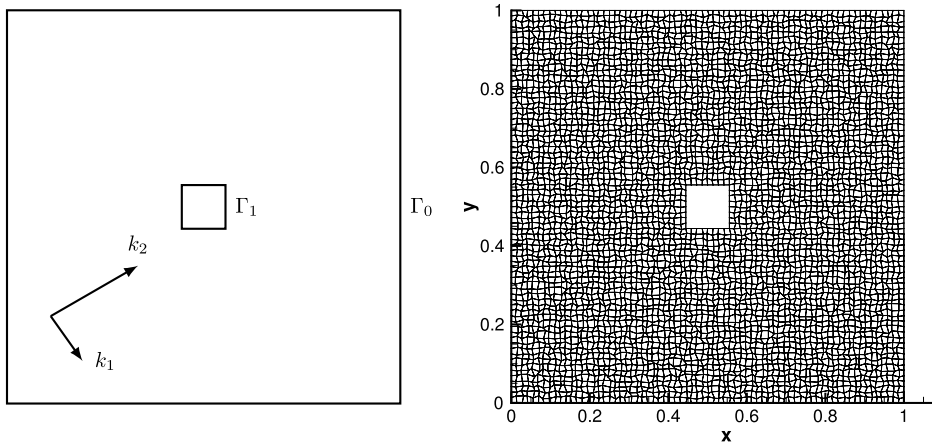


Fig. 7. Left: computational domain; right: distorted random quadrilateral mesh.

outer boundary is referred to as Γ_0 and the internal boundary as Γ_1 . We set the source term $f = 0$, $g_D = 0$ on Γ_0 , $g_D = 2$ on Γ_1 , and take Λ to be the anisotropic tensor

$$\Lambda = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 100 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}, \quad (45)$$

where $\theta = -\pi/6$.

The computational domain and mesh are described in Fig. 7. The numerical solution are shown in the left part of Fig. 8 while the right part is contour. This solution has a sharp gradient near the internal boundary which causes undershoots and overshoots in the numeral solutions. The minimum value is close to zero and the maximum value is 1.98. Hence, our

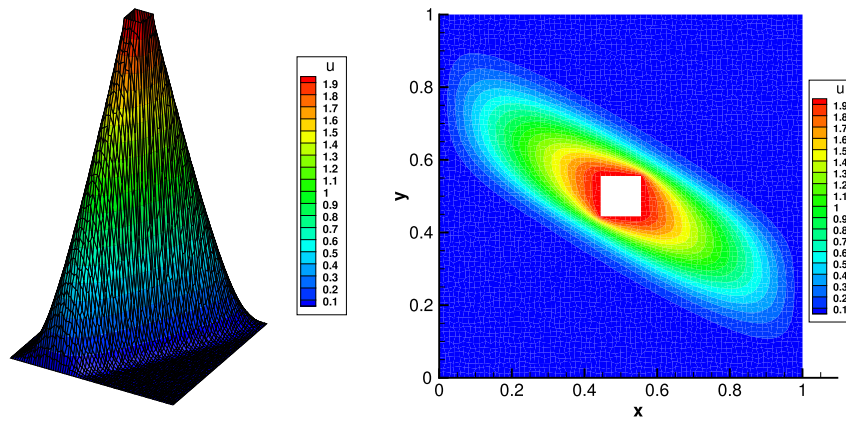


Fig. 8. Left: colormap of numerical solution; right: solution contour.

Table 7

The number of nonlinear iterations on random quadrilateral mesh.

h	VPPS-P	VPPS-AA					
	$m = 1$	$m = 2$	$m = 3$	$m = 5$	$m = 7$	$m = 9$	$m = 11$
1/18	699	591	311	176	138	151	154
1/36	705	611	314	188	201	208	229
1/72	1719	1580	892	410	462	386	319
1/144	311	263	225	162	130	125	227

method produces the non-negative discrete solutions. The number of nonlinear iterations in the Picard and AA methods with various values of m are presented in Table 7, which shows the efficiency of the AA method, however, it will increase the number of iterations if m is too large.

6.5. Spherical nonlinear heat wave

Consider a spherically symmetric heat wave, propagating in a cold medium and governed by the following nonlinear parabolic equation,

$$\frac{\partial T}{\partial t} = \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 \kappa \frac{\partial T}{\partial r}) \quad (46)$$

where T stands for the temperature and depends only on the spherical radius r . The analytic solution is provided by [7] in terms of the radial position of the wave front r_f and the temperature at the center T_c ,

$$T = \begin{cases} T_c \left(1 - \frac{r^2}{r_f^2}\right)^{\frac{1}{2}}, & r \leq r_f, \\ 0, & \text{otherwise,} \end{cases} \quad (47)$$

where

$$T_c = 2^{-\frac{3}{8}} \xi t^{-\frac{3}{8}}, \quad r_f = \xi t^{\frac{1}{8}}, \quad (48)$$

and $\xi = 2^{\frac{7}{8}} / \sqrt{\pi}$ is a dimensionless constant. This problem has been used as a benchmark example by many authors, say, [2,27].

In order to use the present scheme directly, we simply turn to solve the transformation of the above equation

$$u_t - \operatorname{div}(\kappa(u) \nabla u) + \frac{u}{8t} = 0, \quad (49)$$

where $u = T$, $\kappa(u) = T^2$ and $r = \sqrt{x^2 + y^2}$. The computational domain is chosen to be $[0, 1]^2$ and the point source is placed at its bottom-left corner. The boundaries are treated as homogeneous Dirichlet ones.

Randomly distorted quadrilateral meshes are used in this test, which aims to check whether the scheme suffers the numerical heat-barrier issue and how it maintains the spherical symmetry of solution on distorted meshes that do not possess the same symmetry. The simulation starts at $t = 10^{-8}$ and stops at $t = 0.3$. Meanwhile, the time step is a variable

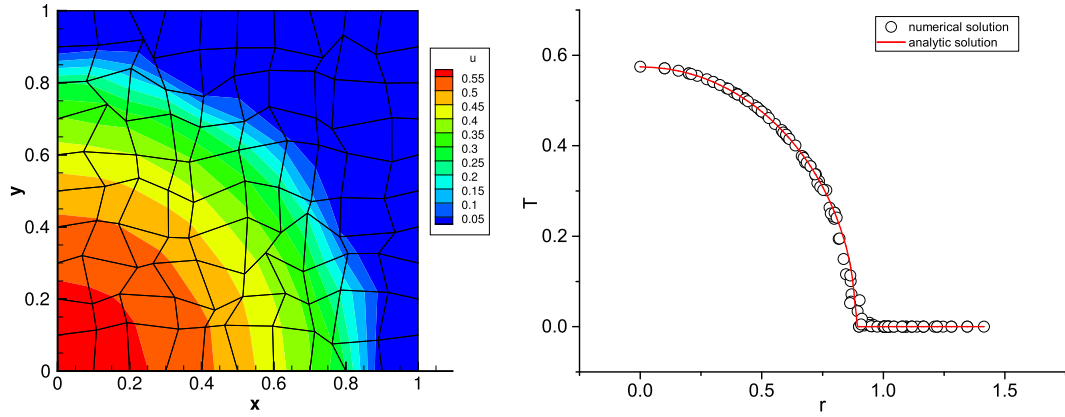


Fig. 9. Results for scheme VPPS on 10×10 random mesh at $t = 0.3$.

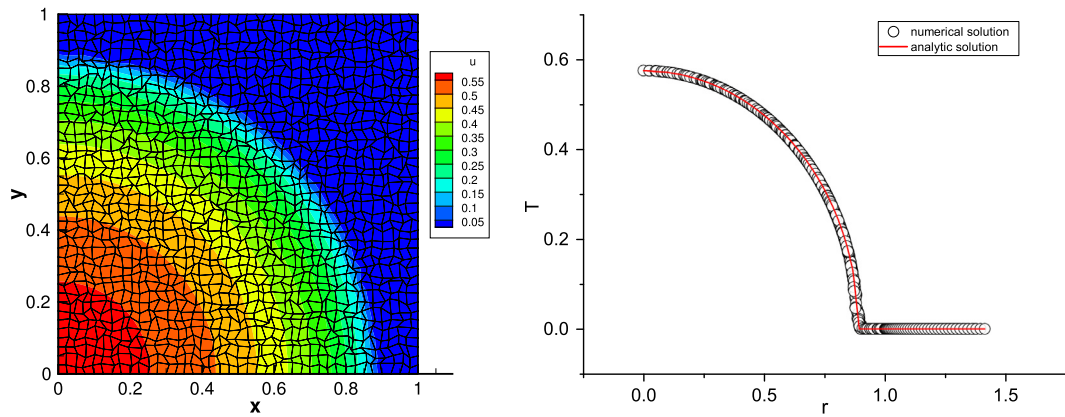


Fig. 10. Results for scheme VPPS on 40×40 random mesh at $t = 0.3$.

one, i.e., it takes $4 \times 10^{-10} \times h^2$ at the beginning and in the following computation it is enlarged (resp. shortened) twenty percent if the number of nonlinear iteration is no more than 5 (resp. greater than 20). The results for VPPS on 10×10 and 40×40 distorted random meshes are shown in Figs. 9 and 10, respectively. The left parts of the figures show the meshes and contours while the right parts are the solution profiles with respect to the spherical radius r . One can see that the new nonlinear FV scheme maintains the spherical symmetry quite well and produces also positive solutions. As for the accuracy, the average convergence rate of the discrete solution on these sequence of meshes is approximately 1.3. The reason for the lower-order convergence is caused by the lower-order regularity of the exact solution.

7. Conclusion

In this paper, a new nonlinear positivity-preserving FV scheme is suggested for the numerical solution of anisotropic diffusion problems on general polygonal grids with star-shaped cells. This scheme differs from existing nonlinear FV schemes in the way that it is of a pure vertex-centered type and relies on a new nonlinear two-point flux approximation that has a fixed stencil. The key ingredient is to abandon the traditional convex decomposition of the co-normal vector, which seems to a unique way to obtain a second-order vertex-centered and positivity-preserving scheme. A simple interpolation algorithm is used for the auxiliary unknowns defined at cell-centers and edge midpoints. Picard method and its Anderson acceleration are discussed for the solution of the nonlinear system. Numerical results show that the present scheme achieves a second-order accuracy, generates nonnegative discrete solutions and moreover, does not suffer the numerical heat-barrier issue. However, if the solution is very close to the machine precision and the mesh distortion is very severe, the truncation error may increase due to the divergence of the Picard iteration. Although this problem can be solved to some extent by choosing a proper parameter m in Anderson acceleration, it is difficult for us to find a choice that works for most cases. The theoretical analysis of nonlinear FV schemes, such as the existence and uniqueness of the discrete solution, the convergence of the fixed-point iteration and the error estimates, is usually very difficult for FV schemes of positivity-preserving type and not touched here. The future works will focus on those topics and the extension of the present scheme to arbitrary polyhedral grids.

Acknowledgements

The authors thank the anonymous reviewers for their careful readings and useful suggestions.

References

- [1] D.G. Anderson, Iterative procedures for nonlinear integral equations, *J. Assoc. Comput. Mach.* 12 (1965) 547–560.
- [2] M.M. Basko, J. Maruhn, A. Tauschwitz, An efficient cell-centered diffusion scheme for quadrilateral grids, *J. Comput. Phys.* 228 (2009) 2175–2193.
- [3] X. Blanc, E. Labourasse, A positive scheme for diffusion problems on deformed meshes, *Z. Angew. Math. Mech.* 96 (2015) 660–680.
- [4] F. Brezzi, K. Lipnikov, V. Simoncini, A family of mimetic finite difference methods on polygonal and polyhedral meshes, *Math. Models Methods Appl. Sci.* 15 (2005) 1533–1551.
- [5] C. Buet, S. Cordier, On the nonexistence of monotone linear schema for some linear parabolic equations, *C. R. Acad. Sci. Paris, Ser. I* 340 (2005) 399–404.
- [6] J. Camier, F. Hermeline, A monotone nonlinear finite volume method for approximating diffusion operators on general meshes, *Int. J. Numer. Methods Eng.* 107 (2016) 496–519.
- [7] Y.B. Zel'dovich, Y.P. Raizer, *Physics of Shock Waves and High-Temperature Hydrodynamic Phenomena*, vol. II, Academic Press, New York, 1967, pp. 864–880.
- [8] J. Droniou, Finite volume schemes for diffusion equations: introduction to and review of modern methods, *Math. Models Methods Appl. Sci.* 24 (2014) 1575–1619.
- [9] R. Eymard, G. Henry, R. Herbin, F. Hubert, R. Klöforn, G. Manzini, 3D benchmark on discretization schemes for anisotropic diffusion problems on general grids, in: J. Fort, J. Furst, J. Halama, R. Herbin, F. Hubert (Eds.), *Finite Volumes for Complex Applications VI*, Springer, 2011, pp. 895–930.
- [10] H. Fang, Y. Saad, Two classes of multisecant methods for nonlinear acceleration, *Numer. Linear Algebra Appl.* 16 (2009) 197–221.
- [11] Z. Gao, J. Wu, A linearity-preserving cell-centered scheme for the heterogeneous and anisotropic diffusion equations on general meshes, *Int. J. Numer. Methods Fluids* 67 (2011) 2157–2183.
- [12] Z. Gao, J. Wu, A small stencil and extremum-preserving scheme for anisotropic diffusion problems on arbitrary 2D and 3D meshes, *J. Comput. Phys.* 250 (2013) 308–331.
- [13] Z. Gao, J. Wu, A second-order positivity-preserving finite volume scheme for diffusion equations on general meshes, *SIAM J. Sci. Comput.* 37 (2015) A420–A438.
- [14] I.V. Kapyrin, A family of monotone methods for the numerical solution of three-dimensional diffusion problems on unstructured tetrahedral meshes, *Dokl. Math.* 76 (2007) 734–738.
- [15] E. Keilegavlen, J. Nordbotten, I. Aavatsmark, Sufficient criteria are necessary for monotone control volume methods, *Appl. Math. Lett.* 22 (2009) 1178–1180.
- [16] K. Lipnikov, M. Shashkov, D. Svyatskiy, Y. Vassilevski, Monotone finite volume schemes for diffusion equations on unstructured triangular and shape-regular polygonal meshes, *J. Comput. Phys.* 227 (2007) 492–512.
- [17] K. Lipnikov, D. Svyatskiy, Y. Vassilevski, Interpolation-free monotone finite volume method for diffusion equations on polygonal meshes, *J. Comput. Phys.* 228 (2009) 703–716.
- [18] K. Lipnikov, D. Svyatskiy, Y. Vassilevski, Anderson acceleration for nonlinear finite volume scheme for advection–diffusion problems, *SIAM J. Sci. Comput.* 35 (2013) A1120–A1136.
- [19] L. Beirão da Veiga, K. Lipnikov, G. Manzini, *The Mimetic Finite Difference Method for Elliptic PDEs*, Springer, New York, 2014.
- [20] K. Lipnikov, G. Manzini, J.D. Moulton, M. Shashkov, The mimetic finite difference method for elliptic and parabolic problems with a staggered discretization of diffusion coefficient, *J. Comput. Phys.* 305 (2016) 111–126.
- [21] J.M. Nordbotten, I. Aavatsmark, G.T. Eigestad, Monotonicity of control volume methods, *Numer. Math.* 106 (2007) 255–288.
- [22] C. Le Potier, Schéma volumes finis monotones pour des opérateurs de diffusion fortement anisotropes sur des maillages de triangles non structurés, *C. R. Acad. Sci. Paris, Ser. I* 341 (2005) 787–792.
- [23] C. Le Potier, A. Mahamane, A nonlinear correction and maximum principle for diffusion operators discretized using hybrid schemes, *C. R. Acad. Sci. Paris, Ser. I* 350 (2012) 101–106.
- [24] M. Schneider, B. Flemisch, R. Helmig, Monotone nonlinear finite-volume method for nonisothermal two-phase two-component flow in porous media, *Int. J. Numer. Methods Fluids* (2016), <http://dx.doi.org/10.1002/flid.4352>, in press.
- [25] Z. Sheng, G. Yuan, The finite volume scheme preserving extremum principle for diffusion equations on polygonal meshes, *J. Comput. Phys.* 230 (2011) 2588–2604.
- [26] Z. Sheng, G. Yuan, An improved monotone finite volume scheme for diffusion equation on polygonal meshes, *J. Comput. Phys.* 231 (2012) 3739–3754.
- [27] C.D. Sijoy, S. Chaturvedi, TRHD: three-temperature radiation-hydrodynamics code with an implicit non-equilibrium radiation transport using a cell-centered monotonic finite volume scheme on unstructured-grids, *Comput. Phys. Commun.* 190 (2015) 98–119.
- [28] A. Toth, C.T. Kelley, Convergence analysis for Anderson acceleration, *SIAM J. Numer. Anal.* 53 (2015) 805–819.
- [29] H.F. Walker, P. Ni, Anderson acceleration for fixed-point iterations, *SIAM J. Numer. Anal.* 49 (2011) 1715–1735.
- [30] T. Washio, C.W. Oosterlee, Krylov subspace acceleration for nonlinear multigrid schemes, *Electron. Trans. Numer. Anal.* 6 (1997) 271–290.
- [31] J. Wu, Vertex-centered linearity-preserving schemes for nonlinear parabolic problems on polygonal grids, *J. Sci. Comput.* 71 (2017) 499–524.
- [32] J. Wu, Z. Gao, Interpolation-based second-order monotone finite volume schemes for anisotropic diffusion equations on general grids, *J. Comput. Phys.* 275 (2014) 569–588.
- [33] J. Wu, Z. Gao, Z. Dai, A vertex-centered linearity-preserving discretization of diffusion problems on polygonal meshes, *Int. J. Numer. Methods Fluids* 81 (2016) 131–150.
- [34] C. Yang, J.C. Meza, B. Lee, L.W. Wang, KSSOLV—a MATLAB toolbox for solving the Kohn–Sham equations, *ACM Trans. Math. Softw.* 36 (2009) 10.
- [35] G. Yuan, Z. Sheng, Monotone finite volume schemes for diffusion equations on polygonal meshes, *J. Comput. Phys.* 227 (2008) 6288–6312.