

On stabilized finite element methods for linear systems of convection–diffusion–reaction equations

Ramon Codina

ETS d'Enginyers de Camins, Canals i Ports, Universitat Politècnica de Catalunya, Jordi Girona 1-3, Edifici C1, 08034 Barcelona, Spain

Received 1 December 1997

Abstract

A stabilized finite element method for solving systems of convection–diffusion–reaction equations is studied in this paper. The method is based on the subgrid scale approach and an algebraic approximation to the subscales. After presenting the formulation of the method, it is analyzed how it behaves under changes of variables, showing that it relies on the law of change of the matrix of stabilization parameters associated to the method. An expression for this matrix is proposed for the case of general coupled systems of equations that is an extension of the expression proposed for a one-dimensional (1D) model problem. Applications of the stabilization technique to the Stokes problem with convection and to the bending of Reissner–Mindlin plates are discussed next. The design of the matrix of stabilization parameters is based on the identification of the stability deficiencies of the standard Galerkin method applied to these two problems. © 2000 Elsevier Science S.A. All rights reserved.

1. Introduction

The purpose of this paper is to study the application of a certain type of stabilized finite element methods to systems of convection–diffusion–reaction equations of the form

$$\mathcal{L}(U) := \frac{\partial}{\partial x_i} (A_i U) - \frac{\partial}{\partial x_i} = F \quad \text{in } \Omega, \quad (1)$$

$$U = 0 \quad \text{on } \partial\Omega, \quad (2)$$

where Ω is the computational domain, U and F are vectors of n_{unk} unknowns and A_i , K_{ij} and S are $n_{\text{unk}} \times n_{\text{unk}}$ matrices ($i, j = 1, \dots, n_{\text{sd}}$). The usual summation convention is implied in (1), with indices running from 1 to the number of space dimensions n_{sd} . We shall refer to the terms on the left-hand-side (LHS) of this equation as the convective, the diffusive and the reactive term. The algebraic bilinear form associated to K_{ij} , $i, j = 1, \dots, n_{\text{sd}}$, is assumed to be positive-definite.

Let $\mathcal{W} := (H_0^1(\Omega))^{n_{\text{unk}}}$. The weak form of the problem consists in finding $U \in \mathcal{W}$ such that

$$a(U, V) - l(V) = 0 \quad \forall V \in \mathcal{W}, \quad (3)$$

where the bilinear form a and the linear form l are defined as

$$a(U, V) := \int_{\Omega} V^t \frac{\partial}{\partial x_i} (A_i U) \, d\Omega + \int_{\Omega} \frac{\partial V^t}{\partial x_i} K_{ij} \frac{\partial U}{\partial x_j} \, d\Omega + \int_{\Omega} V^t S U \, d\Omega, \quad (4)$$

$$l(V) := \int_{\Omega} V^t F \, d\Omega. \quad (5)$$

The Galerkin finite element approximation of this problem is standard. If \mathcal{W}_h is a finite element space to approximate \mathcal{W} , the discrete problem consists in finding $U_h \in \mathcal{W}_h$ such that

$$a(U_h, V_h) - l(V_h) = 0 \quad \forall V_h \in \mathcal{W}_h. \quad (6)$$

It is well known that this formulation lacks stability when the diffusive terms are small, compared either to the convective or to the reactive terms. The purpose of this paper is to analyze certain aspects related to the application of some stabilized finite element techniques to problems (1) and (2). These techniques are described in the following section. They consist in the addition of a residual-based stabilizing term to the basic Galerkin formulation. In particular, we shall concentrate on the subgrid scale method with an algebraic approximation to the subscales, an approach introduced in [1,2]. We present in Section 2 a description of the subgrid scale method in the most general case and also the approximations that lead to the stabilized formulation that will be used in the following, to which we shall refer as the algebraic subgrid scale (ASGS) method. We will also discuss the behavior of the well known SUPG and Galerkin/least-squares (GLS) methods as described for example in [3,4], comparing them with the ASGS formulation.

In Section 3, we study the behavior of the stabilized methods considered under linear changes of variables. It is shown that this simple exercise puts severe restrictions to these methods and to the stabilization parameters on which they depend. In Section 4, a certain expression for these parameters is proposed for different problems. A one-dimensional (1D) convection–diffusion–reaction problem is considered first, obtaining the conditions that the stabilization parameter must verify from the analysis of the discrete maximum principle. The straightforward extension of the expression obtained is then proposed for a general system of equations. Since the resulting stabilization matrix is a matrix function of the coefficients of the differential equation, this approach produces the stabilization parameters of the scalar problem when all the coefficient matrices diagonalize in the same basis, that is, when the problem itself is diagonalizable. However, this general strategy does not work properly for the two problems analyzed next, namely, the Stokes problem with convection and the bending of Reissner–Mindlin plates. The design of the stabilization parameters in these two cases is based on a simple analysis of the lack of stability of the standard Galerkin method when applied to these problems.

In Section 5, a numerical example is presented to check the numerical performance of the stabilized method for Reissner–Mindlin plates and a test is introduced for the general expression of the stabilization matrix. It is shown that the ASGS model, using the numerical parameters proposed in this paper, gives good numerical results for these two problems.

2. Stabilized finite element methods

2.1. The subgrid scale approach

In this section, we present the subgrid scale method in a (slightly) more general version than in the original references [1,2]. Let us split the continuous space \mathcal{W} as $\mathcal{W} = \mathcal{W}_h \oplus \tilde{\mathcal{W}}$, where $\tilde{\mathcal{W}}$ can be in principle any space to complete \mathcal{W}_h in \mathcal{W} . To fix ideas, we may think of $\tilde{\mathcal{W}}$ as the orthogonal complement of \mathcal{W}_h with respect to the L^2 inner product in \mathcal{W} . Since $\tilde{\mathcal{W}}$ represents the component of \mathcal{W} which is not reproduced by the finite element space, we call it the space of subscales or subgrid scales. The continuous equation (3) can now be written as the system

$$a(U_h, V_h) + a(\tilde{U}, V_h) = l(V_h) \quad \forall V_h \in \mathcal{W}_h, \quad (7)$$

$$a(U_h, \tilde{V}) + a(\tilde{U}, \tilde{V}) = l(\tilde{V}) \quad \forall \tilde{V} \in \tilde{\mathcal{W}}, \quad (8)$$

where $U = U_h + \tilde{U}$ and $U_h \in \mathcal{W}_h$, $\tilde{U} \in \tilde{\mathcal{W}}$.

Let n_{el} be the number of elements of the finite element partition of the domain Ω and let Ω^e be the region occupied by the e th element. It is useful for the following to introduce the notation

$$\int_{\Omega'} := \sum_{e=1}^{n_{el}} \int_{\Omega^e}, \quad \int_{\partial\Omega'} := \sum_{e=1}^{n_{el}} \int_{\partial\Omega^e}. \quad (9)$$

Let us assume that the solution of the continuous problem \mathbf{U} is smooth. Integrating by parts within each element domain it is found that problems (7) and (8) can be written as

$$a(\mathbf{U}_h, \mathbf{V}_h) + \int_{\partial\Omega'} \tilde{\mathbf{U}}^t n_i \mathbf{K}_{ij} \frac{\partial \mathbf{V}_h}{\partial x_j} d\Gamma + \int_{\Omega'} \tilde{\mathbf{U}}^t \mathcal{L}^*(\mathbf{V}_h) d\Omega = l(\mathbf{V}_h), \quad (10)$$

$$\int_{\partial\Omega'} \tilde{\mathbf{V}}^t n_i \mathbf{K}_{ij} \frac{\partial}{\partial x_j} (\mathbf{U}_h + \tilde{\mathbf{U}}) d\Gamma + \int_{\Omega'} \tilde{\mathbf{V}}^t \mathcal{L}(\tilde{\mathbf{U}}) d\Omega = \int_{\Omega'} \tilde{\mathbf{V}}^t [\mathbf{F} - \mathcal{L}(\mathbf{U}_h)] d\Omega, \quad (11)$$

where n_i is the i th component of the exterior normal to $\partial\Omega$ and \mathcal{L}^* is the adjoint operator of \mathcal{L} with homogeneous Dirichlet conditions, given below in (20).

Eq. (11) is equivalent to finding $\tilde{\mathbf{U}} \in \tilde{\mathcal{W}}$ such that

$$\mathcal{L}(\tilde{\mathbf{U}}) = \mathbf{F} - \mathcal{L}(\mathbf{U}_h) + \mathbf{V}_{h,\text{ort}} \quad \text{in } \Omega^e, \quad (12)$$

$$\tilde{\mathbf{U}} = \tilde{\mathbf{U}}_{\text{ske}} \quad \text{on } \partial\Omega^e, \quad (13)$$

for $e = 1, \dots, n_{el}$, where $\mathbf{V}_{h,\text{ort}}$ is obtained from the condition that $\tilde{\mathbf{U}}$ must belong to $\tilde{\mathcal{W}}$ (and not to the whole space \mathcal{W}) and $\tilde{\mathbf{U}}_{\text{ske}}$ is a function defined on the element boundaries and such that

$$\mathbf{q}_n := n_i \mathbf{K}_{ij} \frac{\partial}{\partial x_j} (\mathbf{U}_h + \tilde{\mathbf{U}}) \quad (14)$$

is continuous across interelement boundaries, that is to say, the normal component of the fluxes of \mathbf{U} is continuous across these boundaries. Observe that due to this fact the first term in the LHS of (11) vanishes. We call $\tilde{\mathbf{U}}_{\text{ske}}$ the skeleton of $\tilde{\mathbf{U}}$.

Problems (7) and (8) is exactly equivalent to Eqs. (10), (12) and (13). The approximate problem is defined by the way in which problems (12) and (13) is solved as well as by the way in which the functions $\mathbf{V}_{h,\text{ort}}$ and $\tilde{\mathbf{U}}_{\text{ske}}$ are taken. A particular case is described next.

2.2. Algebraic approximation to the subscales

The simplest way to approximate problems (12) and (13) is to take

$$\tilde{\mathbf{U}} \approx \boldsymbol{\tau} [\mathbf{F} - \mathcal{L}(\mathbf{U}_h)] \quad (15)$$

as the solution of this problem, where $\boldsymbol{\tau}$ is a $n_{\text{unk}} \times n_{\text{unk}}$ matrix defined within each element domain that has to be determined. We shall refer to it as the matrix of stabilization parameters. The approximation given by (15) has an implicit assumption on the function $\tilde{\mathbf{U}}_{\text{ske}}$ and the space $\tilde{\mathcal{W}}$, and therefore on the function $\mathbf{V}_{h,\text{ort}}$. In general, $\tilde{\mathbf{U}}$ will be discontinuous across interelement boundaries, so that the fluxes given by (14) will not even be well defined. However, from (10) it is observed that, except for the boundary integral, only the component of $\tilde{\mathbf{U}}$ in $\mathcal{L}(\mathcal{W}_h)$ is needed, where $\mathcal{L}(\mathcal{W}_h)$ is the space of functions of the form $\mathcal{L}(\mathbf{V}_h)$, with $\mathbf{V}_h \in \mathcal{W}_h$. We may think of (15) as the approximation to this component.

To close the approximation, we neglect the interelement boundary terms in (10), so that the problem that has to be solved is finally

$$a(\mathbf{U}_h, \mathbf{V}_h) + \int_{\Omega'} \tilde{\mathbf{U}}^t \mathcal{L}^*(\mathbf{V}_h) d\Omega = l(\mathbf{V}_h), \quad (16)$$

with $\tilde{\mathbf{U}}$ given by (15). With all these assumptions we have arrived to the method proposed in [2] using different arguments. In particular, (15) was derived from an approximation to the Green's function of the problem. This method was also considered in [5] and derived for the scalar diffusion-reaction equation in [6] by using bubble functions.

2.3. The SUPG, GLS and ASGS methods

Let us consider the Galerkin finite element approximation of the problem given by (6). Consider also a stabilized finite element method consisting in adding to the LHS of this equation a term of the form

$$r(\mathbf{U}_h, \mathbf{V}_h) = \int_{\Omega'} \mathcal{P}(\mathbf{V}_h)^t \boldsymbol{\tau} \mathcal{R}(\mathbf{U}_h) d\Omega, \quad (17)$$

where $\mathcal{P}(\mathbf{V}_h)$ is a certain operator applied to the test functions, $\boldsymbol{\tau}$ a matrix of stabilization parameters and $\mathcal{R}(\mathbf{U}_h)$ is the residual of the differential equation, that is to say, $\mathcal{L}(\mathbf{V}_h) - \mathbf{F}$. It is understood that all these terms are computed for each element domain Ω^e .

Most classical stabilization methods for problems (1) and (2) fall within the previous framework, as shown in [7]. For example, for the stationary problem considered in this work, the SUPG, the GLS and the ASGS methods are defined by taking

$$\text{SUPG} : \mathcal{P}(\mathbf{V}_h) = \mathbf{A}_i \frac{\partial \mathbf{V}_h}{\partial x_i}, \quad (18)$$

$$\text{GLS} : \mathcal{P}(\mathbf{V}_h) = \mathcal{L}(\mathbf{V}_h) = \frac{\partial}{\partial x_i} (\mathbf{A}_i \mathbf{V}_h) - \frac{\partial}{\partial x_i} \left(\mathbf{K}_{ij} \frac{\partial \mathbf{V}_h}{\partial x_j} \right) + \mathbf{S} \mathbf{V}_h, \quad (19)$$

$$\text{ASGS} : \mathcal{P}(\mathbf{V}_h) = -\mathcal{L}^*(\mathbf{V}_h) = \mathbf{A}_i^t \frac{\partial \mathbf{V}_h}{\partial x_i} + \frac{\partial}{\partial x_i} \left(\mathbf{K}_{ij}^t \frac{\partial \mathbf{V}_h}{\partial x_j} \right) - \mathbf{S}^t \mathbf{V}_h, \quad (20)$$

where \mathcal{L}^* is the adjoint of the operator \mathcal{L} .

It is important to make some remarks concerning the definition of the SUPG and the GLS methods in (18) and (19). We have taken here the straight extension of these techniques for scalar equations to the vector case. For the SUPG method, the basic idea is to take \mathcal{P} as the convective operator in \mathcal{L} , so as to have control on the convective term of the residual, whereas for the GLS method taking $\mathcal{P} = \mathcal{L}$ leads to a least-squares control of the whole residual. However, it will be shown that these methods do not behave properly when variables are changed. Obviously, this is also true when these methods are applied to the compressible Euler and Navier–Stokes equations as systems of (non-linear) convection–diffusion equations. In the first attempts to apply stabilization techniques to this problem, the misbehavior under changes of variables was overcome by using the transposed coefficient matrices in \mathcal{P} , that is, by taking

$$\text{SUPG} : \mathcal{P}(\mathbf{V}_h) = \mathbf{A}_i^t \frac{\partial \mathbf{V}_h}{\partial x_i}. \quad (21)$$

This was done for example for the Euler equations using conservation variables and the SUPG method in [8], and later in other works such as [9–11], where the use of (21) was taken for granted. In [12], the use of (21) instead of (18) was justified by the fact that it gave ‘superior behavior in non-linear problems of interest’. A later justification was to define the method for the so-called entropy variables [13,14], and then transform to the conservation variables. This change of variables leads again to (21), or to

$$\mathcal{P}(\mathbf{V}) = \mathcal{L}^t(\mathbf{V}) := \mathbf{A}_i^t \frac{\partial \mathbf{V}}{\partial x_i} - \frac{\partial}{\partial x_i} \left(\mathbf{K}_{ij}^t \frac{\partial \mathbf{V}}{\partial x_j} \right) + \mathbf{S}^t \mathbf{V}$$

in the case of the GLS method [15].

3. Change of variables

Symbolically, let us write the stabilized methods as

$$\int_{\Omega} \mathbf{V}^t \mathcal{L}(\mathbf{U}) d\Omega + \int_{\Omega'} \mathcal{P}(\mathbf{V})^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}) d\Omega = \int_{\Omega} \mathbf{V}^t \mathbf{F} d\Omega + \int_{\Omega'} \mathcal{P}(\mathbf{V})^t \boldsymbol{\tau} \mathbf{F} d\Omega. \quad (22)$$

Subscript h has been omitted. It is understood in this equation that the diffusive terms are integrated by parts and the integral of the stabilizing term is evaluated element by element.

We discuss now the conditions under which the stabilized method given by (22) is invariant under changes of variables. Two cases can be distinguished. The first is a ‘true’ change of variables, that is to say, a change of unknowns and force vectors, and the second a change of unknowns only, which is equivalent to a change of variables plus a scaling of the equations.

3.1. Change of unknowns and force vectors

Suppose that

$$\hat{U} = TU, \quad \hat{F} = TF, \quad (23)$$

with T a non-singular matrix of constant coefficients. If we call

$$\hat{\mathcal{L}}(V) = \frac{\partial}{\partial x_i} (TA_i T^{-1} V) - \frac{\partial}{\partial x_i} \left(TK_{ij} T^{-1} \frac{\partial V}{\partial x_j} \right) + TST^{-1} V, \quad (24)$$

then

$$\mathcal{L}(U) = \mathcal{L}(T^{-1} \hat{U}) = T^{-1} \hat{\mathcal{L}}(\hat{U}). \quad (25)$$

The Galerkin contribution to the LHS of (22) can be written as

$$\int_{\Omega} V^t \mathcal{L}(U) d\Omega = \int_{\Omega} V^t T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega = \int_{\Omega} (T^{-1} V)^t \hat{\mathcal{L}}(\hat{U}) d\Omega,$$

from where it follows that the new test function that has to be taken is

$$\hat{V} = T^{-1} V.$$

The stabilizing term in the LHS of (22) can be written now as

$$ST := \int_{\Omega'} \mathcal{P}(V)^t \tau \mathcal{L}(U) d\Omega = \int_{\Omega'} \mathcal{P}(T^t \hat{V})^t \tau T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega. \quad (26)$$

Let us distinguish between the GLS and the ASGS methods. The conclusions obtained below for the GLS method are exactly the same as for the SUPG method, in which case \mathcal{L} is the convective part of \mathcal{L} instead of the whole operator.

• **GLS method:** $\mathcal{P} = \mathcal{L}$

In this case we have that

$$\begin{aligned} ST &= \int_{\Omega} \left[T^{-1} T \mathcal{L}(T^{-1} T T^t \hat{V}) \right]^t \tau T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega \\ &= \int_{\Omega} \hat{\mathcal{L}}(T T^t \hat{V})^t T^{-1} \tau T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega. \end{aligned}$$

From this expression it may be concluded that the GLS method is invariant with respect to changes of variables *only* if

$$T T^t = I, \quad (27)$$

i.e., is invariant to orthogonal changes of variables. But condition (27) is not enough, it is also necessary that τ behaves as

$$\hat{\tau} = T \tau T^t, \quad (28)$$

i.e., τ must be a *matrix function* of the matrix coefficients of the differential equation. If either (27) or (28) do not hold, ‘the’ GLS method with the U variables will be different from ‘the’ GLS method with the \hat{U} variables.

• **ASGS method:** $\mathcal{P} = -\mathcal{L}^*$

From (24) it follows that

$$\hat{\mathcal{L}}^*(V) = -T^{-t} A_i^t T^t \frac{\partial V}{\partial x_i} - \frac{\partial}{\partial x_i} \left(T^{-t} K_{ij}^t T^t \frac{\partial V}{\partial x_j} \right) + T^{-t} S^t T^t V,$$

and therefore the stabilizing term becomes

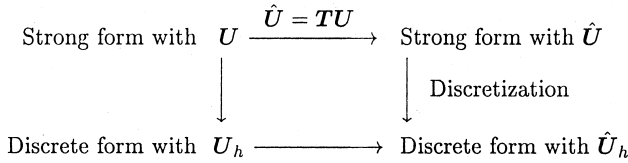
$$\begin{aligned} \text{ST} &= - \int_{\Omega} \left[T^t T^{-t} \mathcal{L}^* \left(T^t \hat{V} \right) \right]^t \tau T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega \\ &= - \int_{\Omega} \hat{\mathcal{L}}^*(\hat{V})^t \tau T^{-1} \hat{\mathcal{L}}(\hat{U}) d\Omega. \end{aligned}$$

From this expression it can be concluded that the ASGS method is invariant to *any* change of variables, provided τ behaves as

$$\hat{\tau} = T \tau T^{-1}, \quad (29)$$

i.e., τ is a *matrix function* of the matrix coefficients.

The situation identified so far can be represented by the following diagram:



For the GLS method the diagram commutes only if T is orthogonal. For the ASGS method it commutes for all T . In both cases, the matrix of stabilizing coefficients τ must be a matrix function of the matrix coefficients.

3.2. Change of unknowns only

Suppose now that

$$\hat{U} = BU, \quad \hat{F} = F, \quad (30)$$

with B a non-singular matrix of constant coefficients. Let

$$\hat{\mathcal{L}}(V) = \frac{\partial}{\partial x_i} (A_i B^{-1} V) - \frac{\partial}{\partial x_i} \left(K_{ij} B^{-1} \frac{\partial V}{\partial x_j} \right) + S B^{-1} V. \quad (31)$$

The Galerkin contribution to the LHS of (22) is

$$\int_{\Omega} V^t \mathcal{L}(U) d\Omega = \int_{\Omega} V^t \hat{\mathcal{L}}(\hat{U}) d\Omega,$$

that is, $\hat{V} = V$. Since the test function does not change, the stabilizing term in the LHS of (22) can be written as

$$\text{ST} = \int_{\Omega} \mathcal{P}(\hat{V})^t \tau \hat{\mathcal{L}}(\hat{U}) d\Omega. \quad (32)$$

Let us distinguish again between the GLS and the ASGS method:

- **GLS method:** $\mathcal{P} = \mathcal{L}$

From (32) we have now that

$$ST = \int_{\Omega} \mathcal{L}(\mathbf{B}^{-1} \mathbf{B} \hat{\mathbf{V}})^t \boldsymbol{\tau} \hat{\mathcal{L}}(\hat{\mathbf{U}}) d\Omega = \int_{\Omega} \hat{\mathcal{L}}(\mathbf{B} \hat{\mathbf{V}})^t \boldsymbol{\tau} \hat{\mathcal{L}}(\hat{\mathbf{U}}) d\Omega.$$

From this equation it follows that for general matrices \mathbf{B} , different sets of variables yield different ‘GLS’ methods.

- **ASGS method:** $\mathcal{P} = -\mathcal{L}^*$

Now we have that

$$\hat{\mathcal{L}}^*(V) = -\mathbf{B}^{-t} \mathbf{A}_i^t \frac{\partial V}{\partial x_i} - \frac{\partial}{\partial x_i} \left(\mathbf{B}^{-t} \mathbf{K}_{ij}^t \frac{\partial V}{\partial x_j} \right) + \mathbf{B}^{-t} \mathbf{S}^t V,$$

and therefore

$$ST = - \int_{\Omega} \left[\mathbf{B}^t \hat{\mathcal{L}}^*(\hat{\mathbf{V}}) \right]^t \boldsymbol{\tau} \hat{\mathcal{L}}(\hat{\mathbf{U}}) d\Omega = - \int_{\Omega} \hat{\mathcal{L}}^*(\hat{\mathbf{V}})^t \mathbf{B} \boldsymbol{\tau} \hat{\mathcal{L}}(\hat{\mathbf{U}}) d\Omega,$$

whereby $\boldsymbol{\tau}$ must verify

$$\hat{\boldsymbol{\tau}} = \mathbf{B} \boldsymbol{\tau}. \quad (33)$$

If (33) does not hold, different variables lead to different methods. Nevertheless, there is the possibility of defining the stabilized method for a particular set of variables \mathbf{U}_{ref} and then to *define* $\boldsymbol{\tau}$ as $\boldsymbol{\tau} = \mathbf{B}_{\text{ref}} \boldsymbol{\tau}_{\text{ref}}$, where $\mathbf{U} = \mathbf{B}_{\text{ref}} \mathbf{U}_{\text{ref}}$.

Changing only the unknowns and not the force vector may be interesting for example if $\mathbf{A}_i \mathbf{B}^{-1}$ are symmetric. This happens for the compressible Navier–Stokes equations (although in this case \mathbf{B} is not a matrix of constant coefficients and the problem is non-linear). The reference unknowns \mathbf{U}_{ref} may be taken as the entropy variables and \mathbf{U} as the conservation variables [14].

If \mathbf{B} is symmetric and positive-definite and $\mathbf{A}_i \mathbf{B}^{-1}$ are symmetric, we may change both unknowns and force vectors and still have symmetric matrices as coefficients for the convective term. To do this, let $\mathbf{B} = \mathbf{L} \mathbf{L}^t$ be the Choleski decomposition of \mathbf{B} . Then

$$\mathbf{L}^{-1} \mathbf{A}_i \mathbf{B} \mathbf{L}^{-t} = \mathbf{L}^{-1} \mathbf{A}_i \mathbf{L}$$

is symmetric, so we may take $\mathbf{T} = \mathbf{L}^{-1}$ in the case of a change of unknowns and force vectors.

A scaling of the differential equation is equivalent to a change of unknowns and force vectors followed by a change of unknowns alone. If \mathbf{T} is the scaling matrix, we may first make a change of unknowns as indicated in (23) with the given scaling matrix and then take $\mathbf{B} = \mathbf{T}^{-1}$ in (30) as matrix for a change of unknowns only. Therefore, the conclusions drawn for the behavior of the GLS and ASGS methods also apply to the scaling of the equations.

4. Applications

In this section we apply the stabilized finite element method for systems of equations discussed above to three different problems of interest. In all the cases, only the matrix of stabilization parameters $\boldsymbol{\tau}$ needs to be defined. The first case corresponds to a general convection–diffusion–reaction system in which the stabilization matrix can be computed from a straightforward extension of the expression for the scalar case derived below. This is possible in particular when all the coefficient matrices of the differential equation diagonalize in the same basis, and therefore the original vector equation can be transformed into a system of uncoupled scalar equations through a change of variables. However, the general expression for $\boldsymbol{\tau}$ obtained from this extension does not work for all the problems of interest. Two examples of this fact are analyzed next, namely, the Stokes problem with convection and the bending of Reissner–Mindlin plates. The design of $\boldsymbol{\tau}$ for these two problems is carried out by looking at the stability problems of the Galerkin

method and trying to improve them. Thus, no general methodology applicable to any convection–diffusion–reaction system of equations will be presented in what follows.

Before considering the three problems mentioned above, it is interesting to analyze what happens when the method is applied to scalar 1D equations. Conditions on τ (now a scalar) can be derived by requiring that the matrix of the final algebraic system be of non-negative type in some limit cases. This leads to an expression of the stabilization parameter that will be used to motivate its counterpart in the problems analyzed next.

4.1. Conditions on τ for the scalar 1D case

Let us consider the 1D problem

$$-k \frac{d^2 u}{dx^2} + a \frac{du}{dx} + su = f, \quad 0 < x < 1, \quad u(0) = u(1) = 0.$$

If the domain is discretized using linear elements of equal length h and the standard Galerkin method is used, the element matrices coming from the diffusive, convective and reactive term are, respectively,

$$A_d^{(e)} = \frac{k}{h} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}, \quad A_c^{(e)} = \frac{a}{2} \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}, \quad A_r^{(e)} = \frac{sh}{6} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}. \quad (34)$$

As in [7], we may obtain a condition for the stability parameter τ by requiring that the matrix of the final algebraic system be of non-negative type, that is, the off-diagonal terms non-positive and the addition of all the coefficients in a row non-negative. This would ensure that the scheme is positive and satisfies the discrete maximum principle. A sufficient condition for this to hold is that the element matrices be of non-negative type. From (34) it follows that this condition is equivalent to

$$-\frac{k}{h} \pm \frac{a}{2} + \frac{sh}{6} \leq 0. \quad (35)$$

When either the SUPG, the GLS or the ASGS methods are used, the final element matrices have the same form as (34) but with modified diffusion, convection and reaction coefficients. The effective parameters using these stabilized methods are

$$\bar{k} = k + \tau a^2, \quad \bar{a} = a - (\xi + 1)\tau a s, \quad \bar{s} = s - \xi \tau s^2,$$

where $\xi = 0$ for the SUPG method, $\xi = -1$ for the GLS method and $\xi = 1$ for the ASGS method. The modified condition (35) now becomes

$$-\frac{\bar{k}}{h} \pm \frac{\bar{a}}{2} + \frac{\bar{s}h}{6} - \tau \left[\frac{\bar{a}^2}{h} \mp \frac{\xi + 1}{2} \bar{a} s + \xi \frac{\bar{s}^2 h}{6} \right] \leq 0. \quad (36)$$

This condition is impossible to fulfill in general, although it provides information about how the different methods behave.

Let us consider the case $s = 0$ first. Condition (36) reduces now to

$$\tau \geq \frac{h}{2a} \left(1 - \frac{1}{Pe} \right), \quad Pe := \frac{ah}{2k} \quad (37)$$

for all the methods.

In the case $a = 0$ condition (36) can be verified only using the ASGS method and provided τ verifies

$$\tau \geq \frac{1}{s} \left(1 - \frac{3}{Ab} \right), \quad Ab := \frac{sh^2}{2k}. \quad (38)$$

Since

$$\begin{aligned}\frac{h}{2a} \left(1 - \frac{1}{Pe}\right) &\leq \frac{h}{2a} \frac{Pe}{Pe+1} = \left(\frac{4k}{h^2} + \frac{2a}{h}\right)^{-1}, \\ \frac{1}{s} \left(1 - \frac{3}{Ab}\right) &\leq \frac{1}{s} \frac{Ab}{Ab+2} = \left(\frac{4k}{h^2} + s\right)^{-1},\end{aligned}$$

we can take

$$\tau = \left(\frac{4k}{h^2} + \frac{2a}{h} + s\right)^{-1}, \quad (39)$$

since this expression verifies the two limiting conditions (37) and (38). Also, it is readily checked that

$$\int_{\Omega^e} \left[N_i + \tau \left(a \frac{dN_i}{dx} - sN_i \right) \right] dx \geq 0,$$

which is needed to keep the sign of f in the components of the discrete force vector.

4.2. Extension to systems

We have found from numerical experiments that the previous expression for τ in (39) yields very good results when extended to scalar equations in multi-dimensional problems, taking in this case a as the Euclidean norm of the velocity \mathbf{a} and h a characteristic length of the element under consideration [7].

Suppose now that the system of equation (1) is diagonalizable, that is, there exists a matrix \mathbf{T} such that all the coefficient matrices in (24) are diagonal. Let

$$\mathbf{K}_0 = (\mathbf{K}_{ij} \mathbf{K}_{ij})^{1/2}, \quad \mathbf{A}_0 = (\mathbf{A}_i \mathbf{A}_i)^{1/2}, \quad \mathbf{S}_0 = (\mathbf{S} \mathbf{S})^{1/2}. \quad (40)$$

Then, the matrix

$$\boldsymbol{\tau} = \left[\frac{c_1}{h^2} \mathbf{K}_0 + \frac{c_2}{h} \mathbf{A}_0 + c_3 \mathbf{S}_0 \right]^{-1} \quad (41)$$

is a matrix function of the coefficient matrices that provides the optimal stabilization parameters for each scalar equation when the system is diagonalized. For linear elements we have found that the constants c_1 and c_2 in (41) may be taken as $c_1 = 4$ and $c_2 = 2$, as for the 1D case. The constant c_3 is taken as 1 in all the cases. For scalar equations it is easy to see that this is mandatory if the instabilities due to dominant reaction terms are to be corrected. A particular example of this fact for systems is the bending of Reissner–Mindlin plates discussed below.

The general expression (41) cannot be applied to an arbitrary system of convection–diffusion–reaction equations. It is obviously effective in the case of diagonalizable systems, and in a numerical example we shall see that it is also useful for systems obtained from the scaling of a diagonal one. However, it does not work for the two very important examples considered in what follows.

4.3. Stokes problem with convection

The first example of the failure of (41) is the generalized Stokes problem

$$-v\Delta \mathbf{u} + \mathbf{a} \cdot \nabla \mathbf{u} + \nabla p = \mathbf{f}, \quad (42)$$

$$\nabla \cdot \mathbf{u} = 0, \quad (43)$$

where \mathbf{u} is the velocity field, p the pressure, \mathbf{f} the vector of body forces, v the kinematic viscosity and \mathbf{a} a given advection velocity, that we assume divergence free to simplify the exposition.

For the sake of simplicity, let us consider the 2D case. Eqs. (42) and (43) can be written as a system of the type (1), the only difference being that the form associated to the diffusion matrices is only positive

semi-definite because there is no diffusion for the pressure, and therefore no boundary conditions can be applied to it. Let $\mathbf{u} = \mathbf{0}$ be the boundary condition for the velocity. The coefficient matrices are now given by

$$\mathbf{K}_{11} = \mathbf{K}_{22} = \begin{bmatrix} v & 0 & 0 \\ 0 & v & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{A}_1 = \begin{bmatrix} a_1 & 0 & 1 \\ 0 & a_1 & 0 \\ 1 & 0 & 0 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} a_2 & 0 & 0 \\ 0 & a_2 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad (44)$$

and $\mathbf{K}_{12} = \mathbf{K}_{21} = \mathbf{S} = \mathbf{0}$.

Numerical experiments indicate that expression (41) to compute the matrix of stabilization parameters does not work properly for this problem. Clearly, this formula for τ does not satisfy condition (33) and therefore the definition of a stabilization method based on it depends on the reference variables chosen and the scaling of the system. If we take $\hat{p} = p/\alpha$ and multiply Eq. (43) by α , the unit coefficients in the \mathbf{A}_i matrices in (44) become α , and the velocity does not change. In other words, if the unit coefficients in the \mathbf{A}_i matrices are multiplied by α , the velocity solution is *exactly the same*. This suggests neglecting these coefficients at the moment of computing τ .

Let us examine which is the lack of stability of the Galerkin method applied to problems (42) and (43). Let $\mathbf{U}_h = [u_{1,h}, u_{2,h}, p_h]^t$ and $\mathbf{V}_h = [v_{1,h}, v_{2,h}, q_h]^t$. The bilinear form associated to the problem given by (4) is now

$$a(\mathbf{U}_h, \mathbf{V}_h) = v \int_{\Omega} \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, d\Omega + \int_{\Omega} (\mathbf{a} \cdot \nabla \mathbf{u}_h) \cdot \mathbf{v}_h \, d\Omega - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} q_h \nabla \cdot \mathbf{u}_h \, d\Omega.$$

Taking $\mathbf{v}_h = \mathbf{u}_h$ and $p_h = q_h$ and denoting by $\|\cdot\|$ the L^2 norm we get

$$a(\mathbf{U}_h, \mathbf{U}_h) = v \|\nabla \mathbf{u}_h\|^2, \quad (45)$$

which determines the stability provided by the Galerkin method. It is observed that the pressure and the convective term are out of control. The stability for the pressure has to be explicitly required by imposing that the finite element spaces to interpolate the velocity and the pressure satisfy the classical inf-sup or Babuška–Brezzi stability condition. To have control on the convective term a sort of streamline diffusion has to be introduced in one way or another.

In this case, the adjoint of the operator \mathcal{L} associated to (42) and (43) is

$$\mathcal{L}^*(\mathbf{V}_h) = \begin{bmatrix} -v\Delta \mathbf{v}_h - \mathbf{a} \cdot \nabla \mathbf{v}_h - \nabla q_h \\ -\nabla \cdot \mathbf{v}_h \end{bmatrix}, \quad (46)$$

and if we take τ as

$$\tau = \text{diag}(\tau_1, \tau_1, \tau_2) \quad (47)$$

the terms to be added to (45) when the ASGS is used are

$$- \int_{\Omega'} \mathcal{L}^*(\mathbf{U}_h)^t \tau \mathcal{L}(\mathbf{U}_h) \, d\Omega = \int_{\Omega'} \left(\tau_1 |\mathbf{a} \cdot \nabla \mathbf{u}_h + \nabla p_h|^2 - \tau_1 v^2 |\Delta \mathbf{u}_h|^2 + \tau_2 |\nabla \cdot \mathbf{u}_h|^2 \right) \, d\Omega. \quad (48)$$

The previous analysis of the 1D model problem and the comments on the extension to multi-dimensional scalar equations suggest to take

$$\tau_1 = \left(\frac{c_1 v}{h^2} + \frac{c_2 |\mathbf{a}|}{h} \right)^{-1}, \quad (49)$$

where c_1 and c_2 are constants. For the same pressure scaling argument as before, we have not considered the pressure in the design of τ_1 . Concerning τ_2 , it helps to improve the control on the divergence of the velocity and is found to be effective in practice [16,17], but for our purposes we can take $\tau_2 = 0$.

The negative sign in the velocity Laplacian term in (48) can be compensated by the control on the velocity gradient given by (45) and using the standard inverse estimate (see e.g., [18]):

$$\|\Delta \mathbf{u}_h\|_{\Omega^e} \leq C_{\text{inv}} \frac{1}{h} \|\nabla \mathbf{u}_h\|_{\Omega^e}. \quad (50)$$

From this and the expression for τ_1 given in (49) we have that

$$-\tau_1 v^2 \|\Delta \mathbf{u}_h\|_{\Omega^e}^2 \geq -\frac{v}{c_1} C_{\text{inv}}^2 \|\nabla \mathbf{u}_h\|_{\Omega^e}^2.$$

We may take $c_1 = 2C_{\text{inv}}^2$ in the definition of τ_1 (in fact, any value $c_1 \geq C_{\text{inv}}^2$ suffices when $|\mathbf{a}| \neq 0$). Obviously, C_{inv} has to be estimated in a general situation. As in the previous cases, we have found that $c_1 = 4$ and $c_2 = 2$ are effective choices for linear elements.

4.4. Reissner–Mindlin plates

The general stabilization method described in the previous sections is applied here to the problem of bending of Reissner–Mindlin plates. In this case, the lack of stability of the standard Galerkin method is demonstrated by the locking effect when the thickness of the plate becomes small.

Stability of Reissner–Mindlin plate elements is only known in few cases and after an important analytical effort (see [19–21] for the analysis of some linear elements). It is shown here that the ASGS method yields a stabilized finite element method for which a simple stability estimate can be obtained. Moreover, this stability analysis dictates how the stability parameter must behave.

Let w be the transverse deflection of the plate and $\boldsymbol{\theta} = [\theta_1, \theta_2]^t$ the rotation vector. Suppose that the plate is clamped. Then, the problem to be solved is

$$-[k_1 \Delta \boldsymbol{\theta} + k_2 \nabla(\nabla \cdot \boldsymbol{\theta})] + \frac{1}{\varepsilon} (\boldsymbol{\theta} - \nabla w) = 0 \quad \text{in } \Omega, \quad (51)$$

$$\frac{1}{\varepsilon} \nabla \cdot (\boldsymbol{\theta} - \nabla w) = q \quad \text{in } \Omega, \quad (52)$$

and $\boldsymbol{\theta} = \mathbf{0}, w = 0$ on $\partial\Omega$. In Eqs. (51) and (52) q is a properly scaled load and

$$k_1 = \frac{E}{24(1+\nu)}, \quad k_2 = \frac{E}{24(1-\nu)}, \quad \varepsilon = \frac{2(1+\nu)}{E\kappa} t^2, \quad (53)$$

and E is the Young modulus, ν the Poisson ratio, κ the shear correction factor and t the plate thickness. We shall write the shear strain as

$$\boldsymbol{\gamma} := \nabla w - \boldsymbol{\theta}. \quad (54)$$

When $\varepsilon \rightarrow 0$ (that is, when $t \rightarrow 0$), the solution w of problems (51) and (52) should converge to the solution of the Lagrange equation of the Kirchhoff plate theory

$$(k_1 + k_2) \Delta \Delta w = q.$$

However, this does not occur when the Galerkin method is used and instead it is found that w tends to zero (or to a wrongly small function) due to the spurious dominance of the shear terms in (51) and (52).

This problem can be recast in the previous framework of systems of convection–diffusion–reaction equations, now with

$$\mathbf{K}_{11} = \begin{bmatrix} k_1 + k_2 & 0 & 0 \\ 0 & k_1 & 0 \\ 0 & 0 & \frac{1}{\varepsilon} \end{bmatrix}, \quad \mathbf{K}_{22} = \begin{bmatrix} k_1 & 0 & 0 \\ 0 & k_1 + k_2 & 0 \\ 0 & 0 & \frac{1}{\varepsilon} \end{bmatrix}, \quad \mathbf{K}_{12} = \mathbf{K}_{21} = \begin{bmatrix} 0 & \frac{k_2}{2} & 0 \\ \frac{k_2}{2} & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 0 & -\frac{1}{\varepsilon} \\ 0 & 0 & 0 \\ \frac{1}{\varepsilon} & 0 & 0 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -\frac{1}{\varepsilon} \\ 0 & \frac{1}{\varepsilon} & 0 \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} \frac{1}{\varepsilon} & 0 & 0 \\ 0 & \frac{1}{\varepsilon} & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Let $\mathbf{U}_h = [\theta_{1,h}, \theta_{2,h}, w_h]^t$ and $\mathbf{V}_h = [\psi_{1,h}, \psi_{2,h}, v_h]^t$ be the trial solution and the test function, respectively, for the finite element approximation of problems (51) and (52) using equal interpolation for rotations and deflection. We assume that this interpolation is made using continuous piecewise polynomials for both fields. The bilinear form associated to the finite element problem can be written as

$$a(\mathbf{U}_h, \mathbf{V}_h) := \int_{\Omega} [k_1 \nabla \boldsymbol{\theta}_h : \nabla \boldsymbol{\psi}_h + k_2 (\nabla \cdot \boldsymbol{\theta}_h)(\nabla \cdot \boldsymbol{\psi}_h)] d\Omega + \frac{1}{\varepsilon} \int_{\Omega} (\boldsymbol{\theta}_h - \nabla w_h) \cdot (\boldsymbol{\psi}_h - \nabla v_h) d\Omega. \quad (55)$$

Taking $\mathbf{V}_h = \mathbf{U}_h$ a simple stability estimate in the L^2 norm is found for the gradients and the divergence of the rotation. However, the coefficients that multiply their norms are k_1 and k_2 , which are negligible compared to the L^2 norm of the shear strain multiplied by $1/\varepsilon$ when $\varepsilon \rightarrow 0$. Thus, rotations are out of control and the shear term dominates the solution.

Let us formulate now the stabilized method. First, observe that if the coefficient matrices \mathbf{A}_1 and \mathbf{A}_2 are multiplied by a given parameter α and the coefficients $(K_{11})_{33}$ and $(K_{22})_{33}$ by α^2 , then the solution in rotations does not change (the new transverse deflection is $\hat{w} = w/\alpha$). Thus, these terms do not need to be included in the design of the matrix of stabilization parameters $\boldsymbol{\tau}$. Also, based on the analysis of the 1D problem in the case in which diffusion and reaction exist, we take

$$\boldsymbol{\tau} = \text{diag}(\tau, \tau, 0), \quad \tau = \left(\frac{c_1 k}{h^2} + c_3 \frac{1}{\varepsilon} \right)^{-1}, \quad (56)$$

where $k := k_1 + k_2$.

Since in this case the operator associated to problems (51) and (52) is self-adjoint, the stabilization term is given by

$$\begin{aligned} - \int_{\Omega'} \mathcal{L}^*(\mathbf{V}_h)^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}_h) d\Omega &= \int_{\Omega'} \tau \left[k_1 \Delta \boldsymbol{\psi}_h + k_2 \nabla (\nabla \cdot \boldsymbol{\psi}_h) - \frac{1}{\varepsilon} (\boldsymbol{\psi}_h - \nabla v_h) \right] \\ &\quad \cdot \left[-k_1 \Delta \boldsymbol{\theta}_h - k_2 \nabla (\nabla \cdot \boldsymbol{\theta}_h) + \frac{1}{\varepsilon} (\boldsymbol{\theta}_h - \nabla w_h) \right] d\Omega. \end{aligned} \quad (57)$$

Let us obtain now a stability estimate for the solution of the stabilized problem. Taking $\mathbf{V}_h = \mathbf{U}_h$ in (55) and (57), adding these two equations up and using Schwarz inequality we get

$$\begin{aligned} a(\mathbf{U}_h, \mathbf{U}_h) - \int_{\Omega'} \mathcal{L}^*(\mathbf{U}_h)^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}_h) d\Omega &\geq \sum_{e=1}^{n_{el}} \left[k_1 \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + k_2 \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \frac{1}{\varepsilon} \|\gamma_h\|_{\Omega^e}^2 - \tau k_1^2 \|\Delta \boldsymbol{\theta}_h\|_{\Omega^e}^2 - \tau k_2^2 \|\nabla (\nabla \cdot \boldsymbol{\theta}_h)\|_{\Omega^e}^2 - \tau \frac{1}{\varepsilon^2} \|\gamma_h\|_{\Omega^e}^2 \right. \\ &\quad \left. - 2\tau \frac{k_1}{\varepsilon} \|\Delta \boldsymbol{\theta}_h\|_{\Omega^e} \|\gamma_h\|_{\Omega^e} - 2\tau \frac{k_2}{\varepsilon} \|\nabla (\nabla \cdot \boldsymbol{\theta}_h)\|_{\Omega^e} \|\gamma_h\|_{\Omega^e} - 2\tau k_1 k_2 \|\Delta \boldsymbol{\theta}_h\|_{\Omega^e} \|\nabla (\nabla \cdot \boldsymbol{\theta}_h)\|_{\Omega^e} \right]. \end{aligned} \quad (58)$$

Using now the inverse estimate (50) and the fact that for any x and y and for any $\lambda > 0$

$$-2xy \geq -\lambda x^2 - \frac{1}{\lambda} y^2,$$

it follows that

$$\begin{aligned}
 a(\mathbf{U}_h, \mathbf{U}_h) - \int_{\Omega'} \mathcal{L}^*(\mathbf{U}_h)^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}_h) d\Omega & \\
 & \geq \sum_{e=1}^{n_{el}} \left[k_1 \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + k_2 \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \frac{1}{\varepsilon} \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 - \tau k_1^2 \frac{C_{inv}^2}{h^2} \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 - \tau k_2^2 \frac{C_{inv}^2}{h^2} \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 - \tau \frac{1}{\varepsilon^2} \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 \right. \\
 & \quad - \tau \frac{k_1}{\varepsilon} \left(\lambda \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \frac{1}{\lambda} \frac{C_{inv}^2}{h^2} \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 \right) - \tau \frac{k_2}{\varepsilon} \left(\lambda \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \frac{1}{\lambda} \frac{C_{inv}^2}{h^2} \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 \right) \\
 & \quad \left. - \tau \left(k_1^2 \frac{C_{inv}^2}{h^2} \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + k_2^2 \frac{C_{inv}^2}{h^2} \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 \right) \right] \\
 & = \sum_{e=1}^{n_{el}} \left[\beta_1 \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \beta_2 \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \beta_3 \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 \right], \tag{59}
 \end{aligned}$$

where

$$\begin{aligned}
 \beta_1 &= \tau \left[c_1 \frac{k_1 k_2}{h^2} + c_1 \frac{k_1^2}{h^2} + c_3 \frac{k_1}{\varepsilon} - 2 \frac{C_{inv}^2}{h^2} k_1^2 - \lambda \frac{k_1}{\varepsilon} \right], \\
 \beta_2 &= \tau \left[c_1 \frac{k_1 k_2}{h^2} + c_1 \frac{k_2^2}{h^2} + c_3 \frac{k_2}{\varepsilon} - 2 \frac{C_{inv}^2}{h^2} k_2^2 - \lambda \frac{k_2}{\varepsilon} \right], \\
 \beta_3 &= \tau \left[c_1 \frac{k}{h^2} \frac{1}{\varepsilon} + c_3 \frac{1}{\varepsilon^2} - \frac{1}{\varepsilon^2} - \frac{1}{\lambda} \frac{C_{inv}^2}{h^2} \frac{k}{\varepsilon} \right].
 \end{aligned}$$

From the expression of β_3 it turns out that the only way to kill the dominance of the shear in the stability estimate (59) is to take $c_3 = 1$, that is, c_3 must be 1. Numerical experiments confirm this fact: if $c_3 > 1$ locking still occurs, whereas if $c_3 < 1$ locking may also occur but since in this case $c_3 - 1 < 0$ the solution may even have the wrong sign!

Taking for example $\lambda = 1/2$ and $c_1 > 2C_{inv}^2$ it follows that there exists a positive constant C , independent of the physical properties and the thickness of the plate, for which

$$\beta_i \geq C k_i, \quad i = 1, 2 \quad \text{and} \quad \beta_3 \geq C \tau \frac{k}{h^2 \varepsilon}.$$

Using this in (59) we finally get the stability estimate

$$a(\mathbf{U}_h, \mathbf{U}_h) - \int_{\Omega'} \mathcal{L}^*(\mathbf{U}_h)^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}_h) d\Omega \geq C \sum_{e=1}^{n_{el}} \left[k_1 \|\nabla \boldsymbol{\theta}_h\|_{\Omega^e}^2 + k_2 \|\nabla \cdot \boldsymbol{\theta}_h\|_{\Omega^e}^2 + \tau \frac{k}{h^2} \frac{1}{\varepsilon} \|\boldsymbol{\gamma}_h\|_{\Omega^e}^2 \right], \tag{60}$$

for a certain constant C .

The important point is the behavior of the coefficient that multiplies the norm of the shear strain. When $\tau = 0$ (Galerkin method) it is $1/\varepsilon$, and thus it tends to infinity as $\varepsilon \rightarrow 0$, whereas now it tends to k/h^2 .

It has been shown that the constant c_3 in (56) must be 1. It remains to define c_1 , which must satisfy the theoretical condition $c_1 > 2C_{inv}^2$. As for the previous problems, we could take $c_1 = 4$ for linear elements. However, the numerical solution is sensitive to this parameter and therefore further analysis is needed in order to determine the optimal value of this constant. An example to quantify this is presented in the following section.

It is interesting to consider the particular case of linear elements or rectangular bilinear elements. The second derivatives within each element are zero and only the shear component in the stabilizing term in (57) remains. Adding this with the original bilinear form in (55) it is found that

$$\begin{aligned}
 a(\mathbf{U}_h, \mathbf{V}_h) - \int_{\Omega'} \mathcal{L}^*(\mathbf{V}_h)^t \boldsymbol{\tau} \mathcal{L}(\mathbf{U}_h) d\Omega &= \int_{\Omega} [k_1 \nabla \boldsymbol{\theta}_h : \nabla \boldsymbol{\psi}_h + k_2 (\nabla \cdot \boldsymbol{\theta}_h) (\nabla \cdot \boldsymbol{\psi}_h)] d\Omega \\
 &\quad + \int_{\Omega'} \left(\frac{1}{\varepsilon} - \tau \frac{1}{\varepsilon^2} \right) (\boldsymbol{\theta}_h - \nabla w_h) \cdot (\boldsymbol{\psi}_h - \nabla v_h) d\Omega.
 \end{aligned}$$

From the expression (56) of τ it turns out that

$$\frac{1}{\varepsilon} - \tau \frac{1}{\varepsilon^2} = \left(\varepsilon + \frac{h^2}{c_1 k} \right)^{-1}, \quad (61)$$

and therefore the stabilized method to which we have arrived consists simply in replacing the factor $1/\varepsilon$ of the shear term by (61), that is, we have recovered the common strategy of using a modified shear correction factor or residual bending flexibility (see for example [22] and references therein). It is important to remark that now we have a variational formulation that justifies this technique and from which its consistency and stability can be established.

5. Numerical examples

In this section we present two numerical examples, one corresponding to the scaling of a diagonal system and the other to a Reissner–Mindlin plate. The behavior of the stabilized formulation presented for the Stokes problem with convection is already well known (except for the design of the stability parameter given in (49), see for example [16]).

5.1. Scaling of a diagonal system

Consider the scalar equation

$$-k\Delta u + \mathbf{a} \cdot \nabla u + su = f \quad \text{in } \Omega = [0, 1]^2,$$

with $\mathbf{a} = A(0.4, 0.7)$, $f = 1$ and the boundary condition $u = 0$ on $\partial\Omega$. Let $\varepsilon = 10^{-5}$ and consider also the following situations:

- (a) $k = 1$, $A = \varepsilon$, $s = \varepsilon$. Solution dominated by diffusion.
- (b) $k = \varepsilon$, $A = 1$, $s = \varepsilon$. Solution dominated by convection.
- (c) $k = \varepsilon$, $A = \varepsilon$, $s = 1$. Solution dominated by reaction.
- (d) $k = \varepsilon$, $A = 1$, $s = 1$. Solution dominated by a combination of convection and reaction.

The idea of the following 2D test with 2 unknowns is to combine these uncoupled basic solutions by adding up their equations, except for the diffusion term, which is always taken as $\mathbf{K}_{11} = \mathbf{K}_{22} = \text{diag}(k_1, k_2)$, $\mathbf{K}_{12} = \mathbf{K}_{21} = \mathbf{0}$. The following two systems are considered:

- (1) First equation obtained from the addition of the equations of cases (a) and (b), second equation that of case (b).
- (2) First equation obtained from the addition of the equations of cases (c) and (d), second equation that of case (d).

For system (1), the numerical solution should be that of cases (a) and (b) for the first and second components, respectively, whereas for system (2) it should be that of cases (c) and (d).

The system of equations to be solved is the original diagonal system composed of two scalar equations scaled by the matrix

$$\mathbf{T} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$

This scaling may be split into a change of unknowns and force vectors as indicated in (23) followed by a change of unknowns only as in (30) with $\mathbf{B} = \mathbf{T}^{-1}$. Since matrix τ given by (41) does not satisfy condition (33), the ASGS model is variable dependent. Therefore, this example serves to test the performance of formula (41) to compute τ for a general coupled system.

Numerical results for this example using a uniform mesh of 20×20 bilinear elements are shown in Figs. 1–6 (first and second components of the vector of unknowns). Results of Fig. 1 have been obtained by

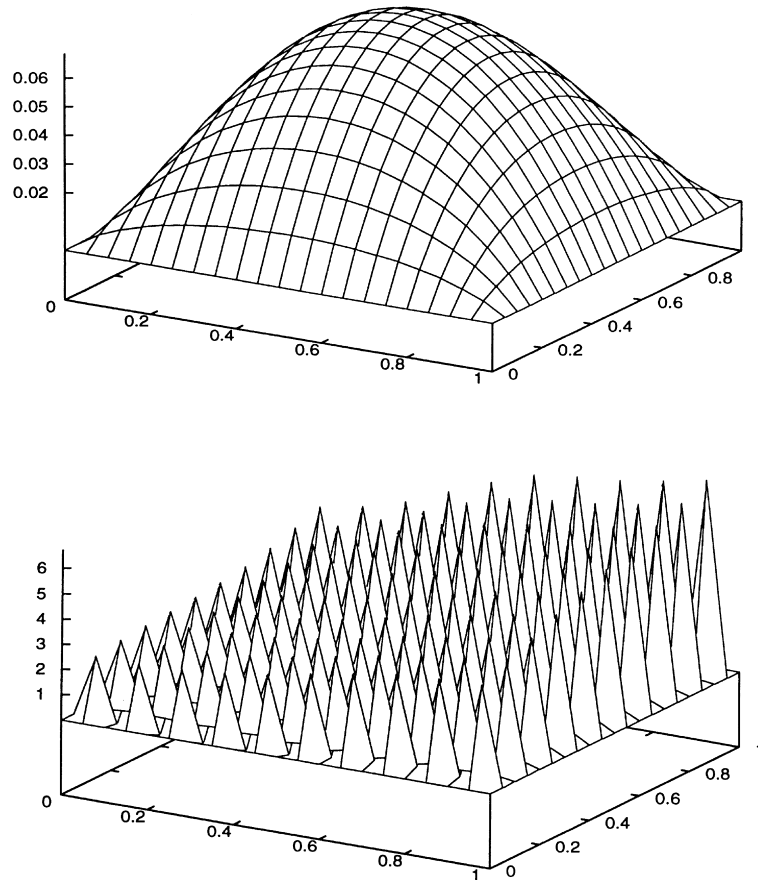


Fig. 1. Results for system 1. ASGS method, τ scalar.

computing τ as $\tau = \tau \mathbf{I}$, where the scalar τ is the minimum of the values obtained for each scalar equation considering only the diagonal terms. This stabilization parameter is not enough for the second unknown and oscillations occur. The solution obtained using τ given by (41) is almost the same as that of the original uncoupled system (Fig. 2), with only small overshoots near the boundary layers. Observe that for this example without reaction terms and using linear elements the ASGS and the GLS methods coincide.

Results of Fig. 3 correspond to system 2 using the GLS method and $\tau = \tau \mathbf{I}$, with τ computed as indicated before. For the first component of the unknown, the oscillations of the standard Galerkin method, that in this case only appear in the neighborhood of the boundary layers, are magnified. However, the second component has been stabilized. Surprisingly, the solution obtained using the matrix form of τ given by (41) is completely oscillatory (Fig. 4). The results obtained using the ASGS method show an improvement with respect to the GLS method when $\tau = \tau \mathbf{I}$ (Fig. 5). In this case, the best solution is obtained with τ computed from (41), in which case only small oscillations for the first unknown are found in the boundary layer where also the second component has overshoots. This is the best one can hope for if the method is not able to recognize that the equations can be in fact uncoupled, that is, if condition (33) is not fulfilled.

In conclusion, results obtained using (41) for this example are very good when using the ASGS method, but not for the GLS method.

5.2. A Reissner–Mindlin plate example

In this example we have solved problems (51) and (52) with $E = 1$, $\nu = 0.2$, $\kappa = 5/6$, $q = 10$ and different values of the plate thickness t . The expression for τ employed is (56), with $c_3 = 1$. The computational domain is again the unit square, discretized now with 200 linear triangles as shown in Fig. 7.

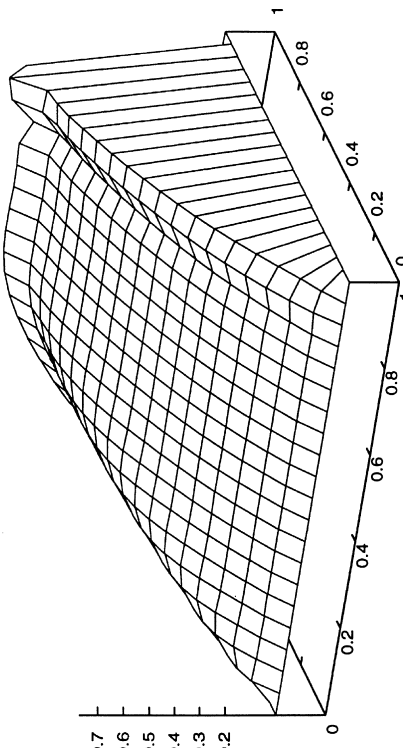
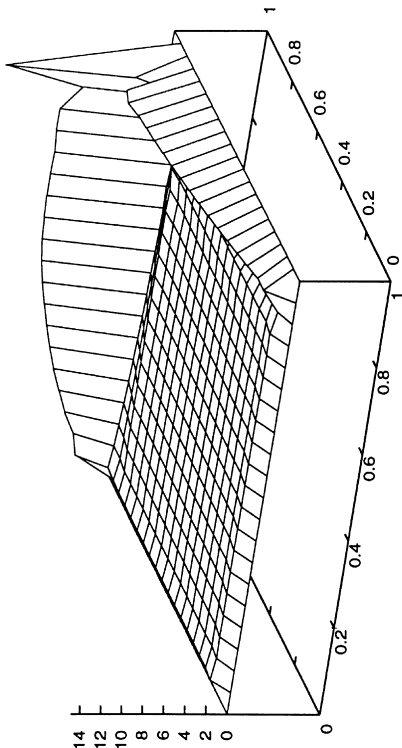


Fig. 3. Results for system 2. GLS method, τ scalar.

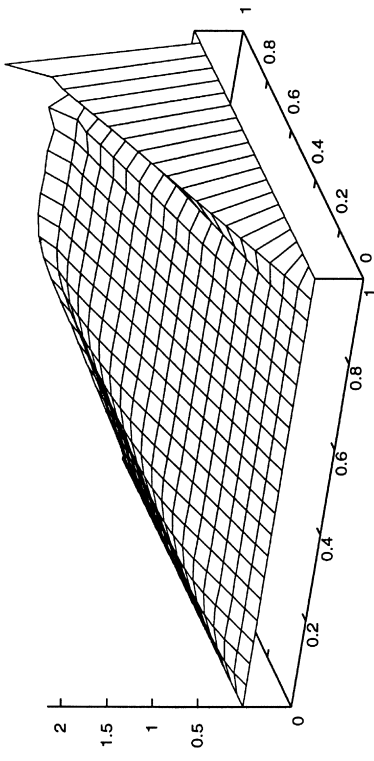
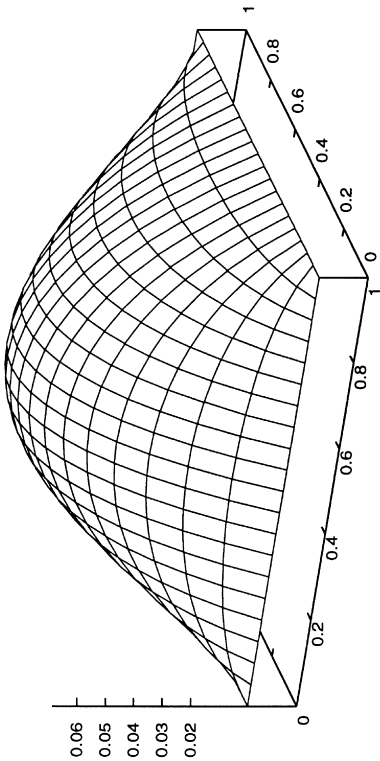


Fig. 2. Results for system 1. ASGS method, τ matrix.

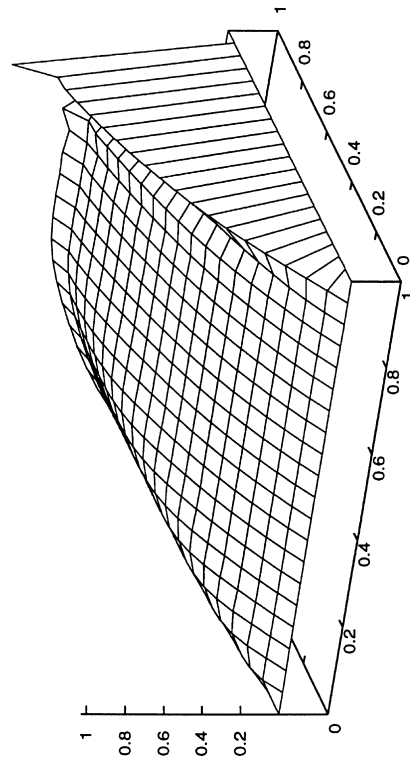
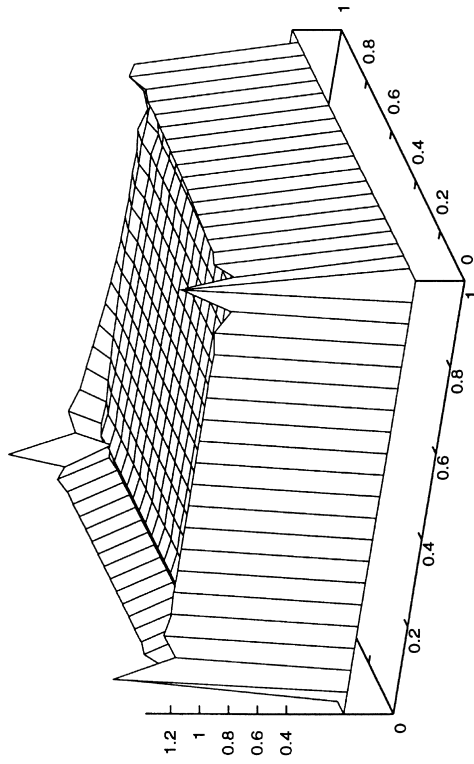


Fig. 5. Results for system 2. ASGS method, τ scalar.

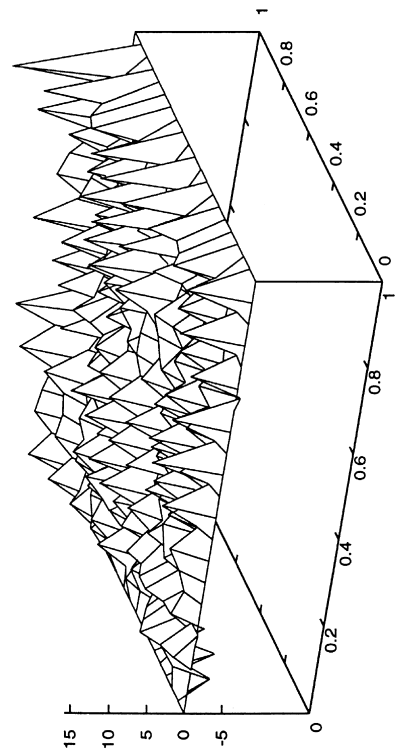
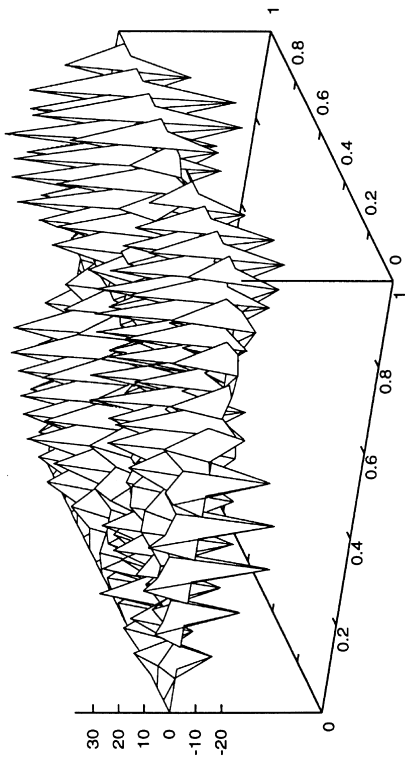


Fig. 4. Results for system 2. GLS method, τ matrix.

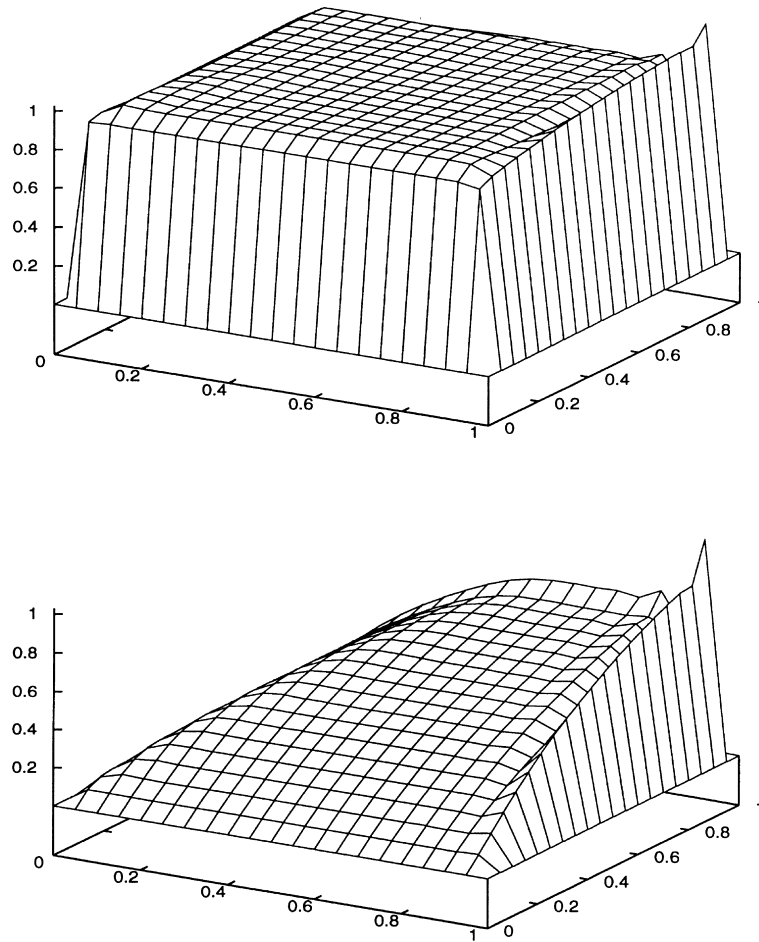


Fig. 6. Results for system 2. ASGS method, τ matrix.

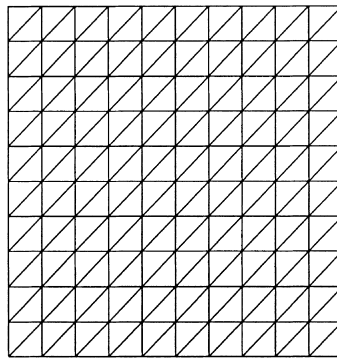


Fig. 7. Mesh for the Reissner–Mindlin plate example.

The central deflection in terms of the plate thickness computed using the standard Galerkin method and the ASGS method with $c_1 = 4$ is shown in Fig. 8. It is observed there how the locking effect occurs using the former, whereas the solution obtained using the latter converges to a non-zero value. The deflection and the second component of the rotation corresponding to this limit case are shown in Fig. 9 (the diagonals of the

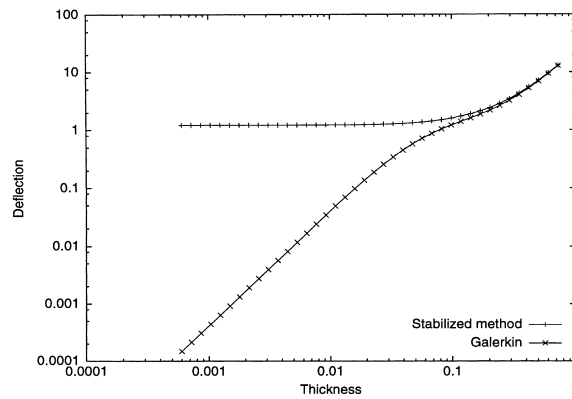


Fig. 8. Results for the Reissner–Mindlin plate. Central deflection vs thickness.

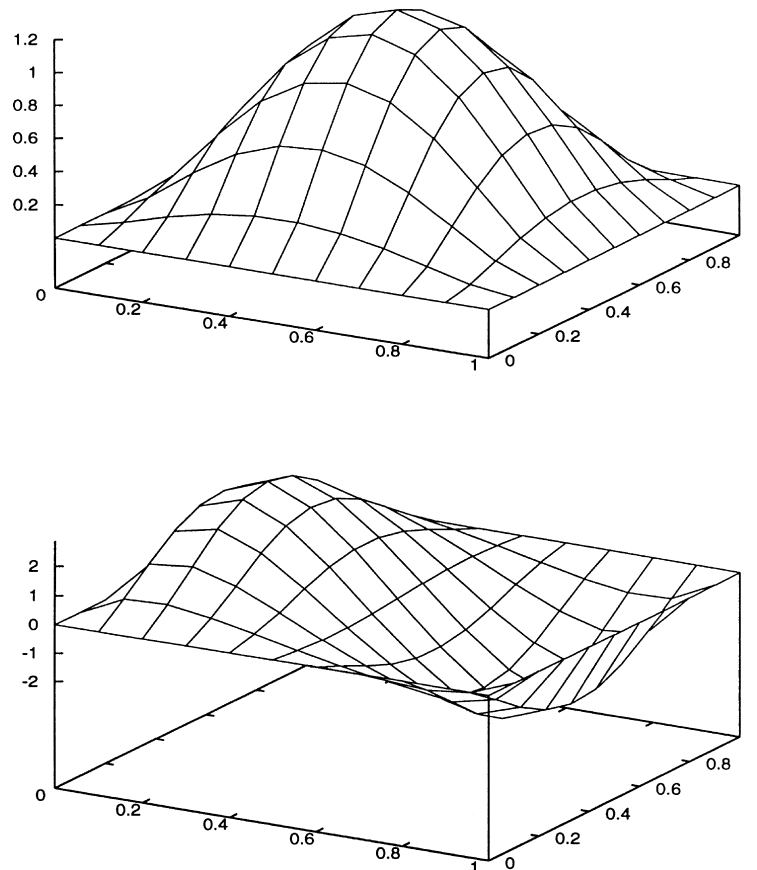


Fig. 9. Results for the Reissner–Mindlin plate for $t = 10^{-4}$. Deflection (top) and second component of the rotation.

triangles are not plotted in these figures). As it was mentioned before, this solution is sensitive to the parameter c_1 in (56). The value of the central deflection in terms of c_1 for $t = 10^{-4}$ is plotted in Fig. 10. In this particular problem, the exact result obtained using the Kirchhoff theory is 1.46, which is obtained for a value of c_1 close to 2.55. For $c_1 = 4$ the central deflection is 1.22. This 16% error is clearly too high for this very simple problem and shows the need for deriving appropriate methods to determine the algorithmic

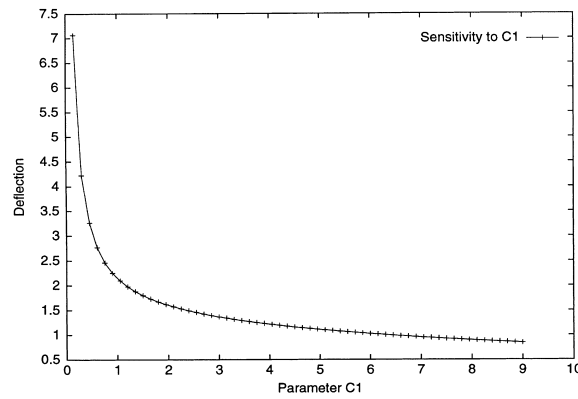


Fig. 10. Results for the Reissner–Mindlin plate. Central deflection vs c_1 .

constants of the formulation or, as it has been shown, the residual bending flexibility. Nevertheless, the method has good stability for all values of h , a feature that the standard Galerkin method lacks.

6. Conclusions

In this paper, several aspects related to the application of a stabilized finite formulation for systems of convection–diffusion–reaction equations have been discussed. First, a general version of the subgrid scale approach has been presented, and particularized next to what has been called algebraic approximation of the subscales. This leads to the ASGS method, which is the stabilized formulation that has been used throughout the paper.

After presenting the relationship between the ASGS and the SUPG and GLS methods, we have discussed the behavior of these methods under changes of variables. This simple exercise has shown that only the ASGS method is invariant under changes of variables, provided that the stabilization matrix τ is a matrix function of the coefficients of the differential equation. The SUPG and the GLS methods are only invariant to orthogonal changes of variables. Also, under changes of variables of the form $\mathbf{U} = \mathbf{B}_{\text{ref}} \mathbf{U}_{\text{ref}}$, keeping constant the force vectors, only the ASGS method allows to design a variable-based method by taking $\tau = \mathbf{B}_{\text{ref}} \tau_{\text{ref}}$.

The most important part of this work is related to the applications of the stabilized formulation to three different problems of interest. Only the matrix of stabilizing parameters needs to be defined for each case. The first conclusion that can be drawn from the analysis of the three problems considered is that no general expression for τ is to be expected. Each problem has its own stability deficiencies that have to be corrected by a proper design of this matrix. The strategy followed here is to identify the terms of the differential equations leading to instability by looking at the coefficients that can be scaled without affecting the numerical results.

Once the terms causing instability have been identified, the design of the stability matrix is based on a simple analysis of what happens for a 1D model problem involving diffusion, convection and reaction. For this case, an expression for the stabilization parameter has been proposed based on the requirement that the matrix of the final algebraic system be of non-negative type in some limit cases, namely, zero reaction and zero convection.

The first problem analyzed is a general system of diagonalizable convection–diffusion–reaction equations. For this problem we have proposed an expression for τ that is the straight extension of what has been found for the 1D case. Since the expression proposed is a matrix function of the matrix coefficients, it provides the stabilization effect of the scalar case for each component of the vector of unknowns. A numerical test introduced here, based on coupling different scalar equations by using linear combinations of them, has shown that the general formula proposed for τ is effective not only for diagonalizable systems, but also when these are scaled.

The next problem analyzed is the Stokes problem with convection, that is, what can be considered as the linear version of the incompressible Navier–Stokes equations. The ASGS method applied to this problem leads to a method similar to the GLS formulation, the only relevant difference being the sign of the viscous term. Also, we have used for this problem an expression for the stabilization parameter motivated by the 1D model problem. It is much simpler than what is commonly used and very effective, both for the numerical analysis and because of the numerical results that it provides (although none of these two aspects has been pursued in this paper).

The final application considered, and perhaps the most innovative of this work, is the bending of Reissner–Mindlin plates using equal interpolation for rotations and normal deflection. For this problem, the ASGS method leads to a non-standard stabilized formulation that is free of locking. We have given here a stability estimate using an expression of the stabilization parameter that, once again, is based on the analysis of the 1D model problem. This estimate is the basic ingredient in the numerical analysis of the method proposed. For linear elements, the stabilized method reduces to the use of a certain shear correction factor or residual bending flexibility, whose use is now justified from a variational standpoint. A numerical example has been presented demonstrating the effectiveness of this method.

Besides the individual conclusions drawn for each problem treated, we want to emphasize that all the methods proposed have been based on the ASGS as a general methodology to derive stabilized formulations when the standard Galerkin method lacks stability. Nevertheless, further improvements may be achieved by better approximating the subscales as solutions of problems (12) and (13), that is to say, by using other subgrid scale models.

References

- [1] T.J.R. Hughes, Multiscale phenomena: Green's function, subgrid scale models, bubbles, and the origins of stabilized formulations, in: M. Morandi Cecchi, K. Morgan, J. Periaux, B.A. Schrefler, O.C. Zienkiewicz (Eds.), *Proceedings of the Ninth International Conference on Finite Elements in Fluids*, Venice, Italy. Dip. di Matematica Pura ed Applicata, Università di Padova, 1995.
- [2] T.J.R. Hughes, Multiscale phenomena: Green's function, the Dirichlet-to-Neumann formulation, subgrid scale models, bubbles and the origins of stabilized formulations, *Comput. Methods Appl. Mech. Engrg.* 127 (1995) 387–401.
- [3] T.J.R. Hughes, M. Mallet, A new finite element method for computational fluid dynamics: III. the generalized streamline operator for multidimensional advective–diffusive systems, *Comput. Methods Appl. Mech. Engrg.* 58 (1986) 305–328.
- [4] T.J.R. Hughes, L.P. Franca, G.M. Hulbert, A new finite element formulation for computational fluid dynamics: VIII. The Galerkin/least-squares method for advective–diffusive equations, *Comput. Methods Appl. Mech. Engrg.* 73 (1989) 173–189.
- [5] L. Franca, S.L. Frey, T.J.R. Hughes, Stabilized finite element methods: I. Application to the advective–diffusive model, *Comput. Methods Appl. Mech. Engrg.* 95 (1992) 253–276.
- [6] L.P. Franca, C. Farhat, Bubble functions prompt unusual stabilized finite element methods, *Comput. Methods Appl. Mech. Engrg.* 123 (1994) 299–308.
- [7] R. Codina, Comparison of some finite element methods for solving the diffusion–convection–reaction equation, *Comput. Methods Appl. Mech. Engrg.* 156 (1998) 185–210.
- [8] T.E. Tezduyar, T.J.R. Hughes, Finite element formulations for convection – dominated flows with particular emphasis on the compressible euler equations, in: *Proceedings of AIAA 21st Aerospace Sciences Meeting*, volume AIAA Paper 83-0125, 1983, 10–13 January 1983/Reno, Nevada.
- [9] S.K. Aliabadi, T.E. Tezduyar, Space-time finite element computation of compressible flows involving moving boundaries and interfaces, *Comput. Methods Appl. Mech. Engrg.* 107 (1993) 209–223.
- [10] G.J. Le Beau, S.E. Ray, S.K. Aliabadi, T.E. Tezduyar, Supg finite element computation of compressible flows with the entropy and conservation variables formulation, *Comput. Methods Appl. Mech. Engrg.* 104 (1993) 397–422.
- [11] A. Soulaïmani, M. Fortin, Finite element solution of compressible viscous flows using conservative variables, *Comput. Methods Appl. Mech. Engrg.* 118 (1994) 319–350.
- [12] T.J.R. Hughes, T.E. Tezduyar, Finite element methods for first - order hyperbolic systems with particular emphasis on the compressible euler equation, *Comput. Methods Appl. Mech. Engrg.* 45 (1984) 217–284.
- [13] C. Johnson, The streamline diffusion finite element method for compressible and incompressible fluid flow, *Von Karman Lecture Series, IV: Computational Fluid Dynamics*, March 1990.
- [14] T.J.R. Hughes, L.P. Franca, M. Mallet, A new finite element formulation for computational fluid dynamics: I. Symmetric forms of the compressible Euler and Navier–Stokes equations and the second law of thermodynamics, *Comput. Methods Appl. Mech. Engrg.* 54 (1986) 223–234.
- [15] G. Hauke, T.J.R. Hughes, A unified approach to compressible and incompressible flows, *Comput. Methods Appl. Mech. Engrg.* 113 (1994) 389–395.
- [16] L.P. Franca, S.L. Frey, Stabilized finite element methods: II. The incompressible Navier–Stokes equations, *Comput. Methods Appl. Mech. Engrg.* 99 (1992) 209–233.

- [17] T.E. Tezduyar, S. Mittal, S.E. Ray, R. Shih, Incompressible flow computations with stabilized bilinear and linear equal-order-interpolation velocity-pressure elements, *Comput. Methods Appl. Mech. Engrg.* 95 (1992) 221–242.
- [18] S.C. Brenner, L.R. Scott, *The Mathematical Theory of Finite Element Methods*, Springer, Germany, 1994.
- [19] D.N. Arnold, R.S. Falk, A uniformly accurate finite element method for the Reissner–Mindlin plate, *SIAM J. Numer. Anal.* 26 (1989) 1276–1290.
- [20] F. Brezzi, M. Fortin, R. Stenberg, Error analysis of mixed interpolated elements for Reissner–Mindlin plates, *Math. Models and Methods in Appl. Sci.* 1 (1991) 125–151.
- [21] L.P. Franca, R. Stenberg, A modification of a low-order Reissner–Mindlin plate bending element, in: J.R. Whiteman (Ed.), *MAFELAP 1990, The Mathematics of Finite Elements and Applications*, Academic Press, New York, 1991, pp. 425–436.
- [22] A. Tessler, T.J.R. Hughes, A three-node Mindlin plate element with improved transverse shear, *Comput. Methods Appl. Mech. Engrg.* 50 (1985) 71–101.