

## CONDITIONING OF FINITE ELEMENT EQUATIONS WITH ARBITRARY ANISOTROPIC MESHES

LENNARD KAMENSKI, WEIZHANG HUANG, AND HONGGUO XU

**ABSTRACT.** Bounds are developed for the condition number of the linear finite element equations of an anisotropic diffusion problem with arbitrary meshes. They depend on three factors. The first factor is proportional to a power of the number of mesh elements and represents the condition number of the linear finite element equations for the Laplacian operator on a uniform mesh. The other two factors arise from the mesh nonuniformity viewed in the Euclidean metric and in the metric defined by the diffusion matrix. The new bounds reveal that the conditioning of the finite element equations with adaptive anisotropic meshes is much better than what is commonly assumed. Diagonal scaling for the linear system and its effects on the conditioning are also studied. It is shown that the Jacobi preconditioning, which is an optimal diagonal scaling for a symmetric positive definite sparse matrix, can eliminate the effects of mesh nonuniformity viewed in the Euclidean metric and reduce those effects of the mesh viewed in the metric defined by the diffusion matrix. Tight bounds on the extreme eigenvalues of the stiffness and mass matrices are obtained. Numerical examples are given.

### 1. INTRODUCTION

It has been amply demonstrated that significant improvements in accuracy can be gained when an appropriately chosen anisotropic mesh is used for the numerical solution of problems exhibiting anisotropic features. However, there exists a general concern in the scientific computing community that an anisotropic mesh, which can contain elements of large aspect ratio, may lead to ill-conditioned linear systems and this could outweigh the accuracy and efficiency improvements gained by anisotropic mesh adaptation. For isotropic mesh adaptation, Bank and Scott [2] (also see Brenner and Scott [3]) show that after proper diagonal scaling, the condition number of finite element equations with an adaptive mesh is essentially the same as that for a uniform mesh. Unfortunately, this result does not apply to anisotropic meshes nor to problems with anisotropic diffusion.

For problems with anisotropic diffusion and arbitrary meshes, several estimates have been developed for the extreme eigenvalues of the stiffness matrix. For example, Fried [7] shows that the largest eigenvalue of the stiffness matrix is bounded by the largest eigenvalues of element stiffness matrices. Shewchuk [18] obtains sharp

---

Received by the editor January 17, 2012 and, in revised form, September 8, 2012, September 17, 2012 and January 4, 2013.

2010 *Mathematics Subject Classification.* Primary 65N30, 65N50, 65F35, 65F15.

*Key words and phrases.* Mesh adaptation, anisotropic mesh, finite element, mass matrix, stiffness matrix, conditioning, extreme eigenvalues, preconditioning, diagonal scaling.

This work was supported in part by the DFG (Germany) under grants KA 3215/1-1 and KA 3215/2-1 and the National Science Foundation (U.S.A.) under grants DMS-0712935 and DMS-1115118.

©2014 American Mathematical Society  
Reverts to public domain 28 years from publication

bounds on the largest eigenvalues of element stiffness matrices for linear triangular and tetrahedral finite elements. More recently, Du et al. [5] developed a bound that can be viewed as a generalization of Shewchuk's result to general dimensions and simplicial finite elements.

Estimation of the smallest eigenvalue for the general case appears to be more challenging. Standard estimates (e.g., see Ern and Guermond [6]) are linearly proportional to the volume of the smallest mesh element, which is typically too pessimistic for nonuniform meshes. Moreover, Apel [1, Sect. 4.3.3] shows that the order of the smallest eigenvalue of the stiffness matrix for a specific, specially designed anisotropic mesh is the same as for a uniform mesh. As a matter of fact, adaptive meshes based on the coefficients of partial differential equations (PDEs) can even improve the conditioning for PDEs with anisotropic diffusion coefficients, as observed by D'Azevedo et al. [4] and Shewchuk [18, Sect. 3.2]. A noticeable approach for obtaining sharper bounds for the smallest eigenvalue is proposed by Fried [7]. The approach employs a continuous generalized eigenvalue problem with an auxiliary density function and its key is to find a lower bound for the smallest eigenvalue of the continuous problem. Bounds for the smallest eigenvalue of the stiffness matrix obtained with Fried's approach are valid for general meshes in any dimension but in  $d \geq 3$  dimensions they are less sharp than those obtained in this paper.

The objective of this paper is threefold. First, we develop tight bounds on the extreme eigenvalues and the condition number of the stiffness matrix for a general diffusion problem with an arbitrary anisotropic mesh. No assumption on the shape or size of mesh elements is made in the development. Our upper bound on the largest eigenvalue can also be expressed in terms of mesh nonuniformity viewed in the metric tensor defined by the diffusion matrix (which will hereafter be referred to as the *mesh  $\mathbb{D}$ -nonuniformity*). It is comparable to those of Shewchuk [18] and Du et al. [5] but is expressed as a sum of patchwise terms instead of elementwise terms as in the aforementioned references. The patchwise nature gives a sharper bound and makes it more convenient to use in the development of diagonal scaling preconditioners. To obtain lower bounds for the smallest eigenvalue of the stiffness matrix we extend Bank and Scott's result [2] to arbitrary meshes. This generalization is not trivial and special effort has to be made to deal with the arbitrariness of the mesh. Along the way we establish anisotropic upper and lower bounds on the extreme eigenvalues of the mass matrix which are much tighter than estimates available in the literature.

The second objective of the paper is to provide a clear geometric interpretation for the obtained bounds on the condition number of the stiffness matrix. These bounds are shown to depend on three factors. The first factor is proportional to a power of the number of mesh elements and represents the condition number of the stiffness matrix for the linear finite element approximation of the Laplacian operator on a uniform mesh. The other two factors arise from the mesh nonuniformity in volume measured in the Euclidean metric (which will be referred to as the *mesh volume-nonuniformity*) and from the mesh  $\mathbb{D}$ -nonuniformity.

The third objective is to study diagonal scaling for the finite element linear system and its effects on the conditioning. We focus on the scaling with the diagonal entries of the matrix (Jacobi preconditioning) since it is an optimal diagonal scaling

for a symmetric positive definite sparse matrix [10, Corollary 7.6 and the following]. We show that the Jacobi preconditioning can eliminate the effects of the mesh volume-nonuniformity and improve those caused by the mesh  $\mathbb{D}$ -nonuniformity, thus significantly reducing the effects of the mesh irregularity on the conditioning. From the practical point of view, this result indicates that a simple diagonal preconditioning can effectively transform the stiffness matrix into a matrix which has a comparable condition number as the one with a uniform mesh.

The outline of the paper is as follows. Section 2 briefly describes a linear finite element discretization of a general anisotropic diffusion problem. Estimation of the extreme eigenvalues and the condition number of the mass matrix is given in Section 3. Section 4 deals with the estimation of the largest eigenvalue of the stiffness matrix. Bounds on the smallest eigenvalue and the condition number of the stiffness matrix and the effects of diagonal scaling are investigated in Section 5. A selection of examples in Section 6 provides a numerical validation for the theoretical findings. Finally, conclusions and further remarks are given in Section 7.

2. LINEAR FINITE ELEMENT APPROXIMATION

We consider the boundary value problem (BVP) of a general diffusion differential equation in the form

$$(1) \quad \begin{cases} -\nabla \cdot (\mathbb{D}\nabla u) = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial\Omega, \end{cases}$$

where  $\Omega$  is a simply connected polygonal or polyhedral domain in  $\mathbb{R}^d$  ( $d \geq 1$ ) and  $\mathbb{D} = \mathbb{D}(\mathbf{x})$  is the diffusion matrix. We assume that  $\mathbb{D}$  is symmetric and positive definite and there exist two positive constants,  $d_{\min}$  and  $d_{\max}$ , such that

$$(2) \quad d_{\min}I \leq \mathbb{D}(\mathbf{x}) \leq d_{\max}I, \quad \forall \mathbf{x} \in \Omega,$$

where the less-than-or-equal sign means that the difference between the right-hand side and left-hand side terms is positive semidefinite.

We are interested in the linear finite element solution of BVP (1). Assume that an affine family  $\{\mathcal{T}_h\}$  of simplicial decompositions of  $\Omega$  is given and denote the associated linear finite element space by  $V^h \subset H_0^1(\Omega)$ . A linear finite element solution  $u_h \in V^h$  to BVP (1) is defined by

$$\int_{\Omega} \nabla v_h \cdot \mathbb{D}\nabla u_h \, d\mathbf{x} = \int_{\Omega} f v_h \, d\mathbf{x}, \quad \forall v_h \in V^h$$

or

$$\sum_{K \in \mathcal{T}_h} \int_K \nabla v_h \cdot \mathbb{D}\nabla u_h \, d\mathbf{x} = \sum_{K \in \mathcal{T}_h} \int_K f v_h \, d\mathbf{x}, \quad \forall v_h \in V^h.$$

Since both  $\nabla u_h$  and  $\nabla v_h$  are constant on  $K$ , we can rewrite the above equation as

$$(3) \quad \sum_{K \in \mathcal{T}_h} |K| \nabla v_h \cdot \mathbb{D}_K \nabla u_h = \sum_{K \in \mathcal{T}_h} \int_K f v_h \, d\mathbf{x}, \quad \forall v_h \in V^h,$$

where  $\mathbb{D}_K$  is the integral average of  $\mathbb{D}$  over  $K$ , i.e.,

$$(4) \quad \mathbb{D}_K = \frac{1}{|K|} \int_K \mathbb{D}(\mathbf{x}) \, d\mathbf{x}.$$

In practice, the integrals in (3) and (4) have to be approximated numerically via a quadrature rule. Although this will change the definition of  $\mathbb{D}_K$  and the right-hand

side term of (3) slightly, the procedure and the results in this paper will remain valid for this situation.

Finite element equation (3) can be expressed in a matrix form. Denoting the numbers of elements and interior vertices of  $\mathcal{T}_h$  by  $N$  and  $N_{vi}$  and assuming that the vertices are ordered in such a way that the first  $N_{vi}$  vertices are the interior vertices, we have

$$(5) \quad \begin{aligned} V^h &= \text{span}\{\phi_1, \dots, \phi_{N_{vi}}\}, \\ u_h &= \sum_j u_j \phi_j, \end{aligned}$$

where  $\phi_j$  is the linear basis function associated with the  $j^{\text{th}}$  vertex. In (5) and hereafter, we use the sum  $\sum_j$  with the index  $j$  ranging over all interior vertices, i.e.,

$$\sum_j = \sum_{j=1}^{N_{vi}}$$

Substituting (5) into (3) and taking  $v_h = \phi_i$  ( $i = 1, \dots, N_{vi}$ ), we obtain the linear algebraic system

$$A\mathbf{u} = \mathbf{f},$$

where  $\mathbf{u} = [u_1, \dots, u_{N_{vi}}]^T$  and the stiffness matrix  $A$  and the right-hand side term  $\mathbf{f}$  are given by

$$(6) \quad \begin{aligned} A_{ij} &= \sum_{K \in \mathcal{T}_h} |K| \nabla \phi_i|_K \cdot \mathbb{D}_K \nabla \phi_j|_K, & i, j &= 1, \dots, N_{vi}, \\ f_i &= \sum_{K \in \mathcal{T}_h} \int_K f \phi_i \, d\mathbf{x}, & i &= 1, \dots, N_{vi}, \end{aligned}$$

and  $\nabla \phi_i|_K$  and  $\nabla \phi_j|_K$  denote the restriction of  $\nabla \phi_i$  and  $\nabla \phi_j$  on  $K$ . Our main goal is to estimate the condition number of stiffness matrix  $A$ .

### 3. MASS MATRIX

To start with, we consider the element mass matrix  $\hat{B}$  for the reference element  $\hat{K}$  (which is assumed to be unitary, i.e.,  $|\hat{K}| = 1$ ),

$$\hat{B} = (\hat{B}_{ij}), \quad \hat{B}_{ij} = \int_{\hat{K}} \hat{\phi}_i \hat{\phi}_j \, d\boldsymbol{\xi} = \frac{1 + \delta_{ij}}{(d+1)(d+2)}, \quad i, j = 1, \dots, d+1,$$

where  $\hat{\phi}_i$ 's are the linear basis functions associated with the vertices of  $\hat{K}$  and  $\delta_{ij}$  is the Kronecker delta function. Matrix  $\hat{B}$  is symmetric and positive definite. Moreover, by direct calculation it can be found that the eigenvalues of  $\hat{B}$  are  $\lambda_1 = \frac{1}{d+1}$  and  $\lambda_2 = \dots = \lambda_{d+1} = \frac{1}{(d+1)(d+2)}$ . Thus,

$$\frac{1}{(d+1)(d+2)} I \leq \hat{B} \leq \frac{1}{d+1} I.$$

**3.1. Condition number of the mass matrix.** Consider the global mass matrix

$$B = (B_{ij}), \quad B_{ij} = \int_{\Omega} \phi_i \phi_j \, d\mathbf{x}, \quad i, j = 1, \dots, N_{vi}.$$

Notice that

$$(7) \quad B_{jj} = \int_{\Omega} \phi_j^2 \, d\mathbf{x} = \sum_{K \in \omega_j} \int_K \phi_j^2 \, d\mathbf{x} = \sum_{K \in \omega_j} \frac{2|K|}{(d+1)(d+2)} = \frac{2|\omega_j|}{(d+1)(d+2)},$$

where  $\omega_j$  is the element patch associated with the  $j^{\text{th}}$  vertex and  $|\omega_j|$  is its volume.

The following theorem gives lower and upper bounds on the condition number of the mass matrix for any dimension and any mesh.

**Theorem 3.1** (Condition number of the mass matrix). *The condition number of the mass matrix for the linear finite elements on a simplicial mesh is bounded by*

$$(8) \quad \frac{\max_j B_{jj}}{\min_j B_{jj}} \leq \kappa(B) \leq (d+2) \frac{\max_j B_{jj}}{\min_j B_{jj}}.$$

*Proof.* For an element  $K$ , let  $\mathbf{u}_K$  be the restriction of the vector  $\mathbf{u}$  on  $K$  and  $B_K$  the element mass matrix. Then,

$$\mathbf{u}^T B \mathbf{u} = \sum_{K \in \mathcal{T}_h} \mathbf{u}_K^T B_K \mathbf{u}_K = \sum_{K \in \mathcal{T}_h} |K| \mathbf{u}_K^T \hat{B} \mathbf{u}_K \leq \frac{1}{d+1} \sum_{K \in \mathcal{T}_h} |K| \|\mathbf{u}_K\|_2^2.$$

Rearranging the sum on the right-hand side according to the vertices and using (7),

$$\mathbf{u}^T B \mathbf{u} \leq \frac{1}{d+1} \sum_{K \in \mathcal{T}_h} |K| \|\mathbf{u}_K\|_2^2 = \frac{1}{d+1} \sum_j u_j^2 |\omega_j| = \frac{d+2}{2} \sum_j u_j^2 B_{jj},$$

which implies

$$\lambda_{\max}(B) \leq \frac{d+2}{2} \max_j B_{jj}.$$

Similarly, we have

$$\mathbf{u}^T B \mathbf{u} \geq \frac{1}{2} \sum_j u_j^2 B_{jj} \quad \text{and} \quad \lambda_{\min}(B) \geq \frac{1}{2} \min_j B_{jj}.$$

Moreover, it is easy to show that

$$\lambda_{\max}(B) \geq \max_j B_{jj} \quad \text{and} \quad \lambda_{\min}(B) \leq \min_j B_{jj}.$$

Combining the above estimates gives

$$\begin{aligned} \max_j B_{jj} \leq \lambda_{\max}(B) &\leq \frac{d+2}{2} \max_j B_{jj}, \\ \frac{1}{2} \min_j B_{jj} \leq \lambda_{\min}(B) &\leq \min_j B_{jj}, \end{aligned}$$

from which the estimate (9) follows. □

From (7), Theorem 3.1 implies

$$(9) \quad \frac{|\omega_{\max}|}{|\omega_{\min}|} \leq \kappa(B) \leq (d+2) \frac{|\omega_{\max}|}{|\omega_{\min}|},$$

where  $|\omega_{\max}| = \max_j |\omega_j|$  and  $|\omega_{\min}| = \min_j |\omega_j|$ . If the mesh is almost uniform, then  $|\omega_{\max}|/|\omega_{\min}| = \mathcal{O}(1)$  and  $\kappa(B) = \mathcal{O}(1)$ . On the other hand,  $|\omega_{\max}|/|\omega_{\min}|$  and  $\kappa(B)$  can become large for nonuniform meshes.

**3.2. Relation to the estimates in the literature.** If we denote the maximum number of mesh elements in a patch by  $p_{\max}$ , then

$$|K_{\min}| \leq |\omega_j| \leq p_{\max} |K_{\max}|, \quad \forall j = 1, \dots, N_{vi}$$

and estimate (9) implies

$$(10) \quad \kappa(B) \leq (d + 2)p_{\max} \frac{|K_{\max}|}{|K_{\min}|},$$

which is the bound obtained by Fried [7, inequality (24)].

Moreover, for an isotropic mesh,

$$|K_{\max}| \propto h_{\max}^d, \quad |K_{\min}| \propto h_{\min}^d,$$

where  $h_{\max}$  and  $h_{\min}$  are the largest and smallest element diameters. Substituting this into (10) gives

$$(11) \quad \kappa(B) \leq C \left( \frac{h_{\max}}{h_{\min}} \right)^d,$$

which is precisely the standard estimate found in the literature (e.g., [6, Rem. 9.10]).

For anisotropic meshes, on the other hand, the new estimate (9) is much tighter than both Fried’s estimate (10) and the standard estimate (11), since large  $|K_{\max}|/|K_{\min}|$  and  $h_{\max}/h_{\min}$  do not necessarily imply large  $|\omega_{\max}|/|\omega_{\min}|$  for those meshes.<sup>1</sup> Furthermore, estimate (9) also provides a tight lower bound, which is not available with (10) and (11).

**3.3. Diagonal scaling for the mass matrix.** It is known [10, Corollary 7.6 and the following] that for a symmetric positive definite sparse matrix, scaling by its diagonal entries (Jacobi preconditioning) is an optimal diagonal preconditioning (up to a constant depending on the maximum number of nonzeros per column and row of the matrix). We are interested in a bound on the condition number after such preconditioning.

For a diagonal scaling  $S = (s_j)$ , similarly to Theorem 3.1 we obtain

$$\frac{\max_j s_j^{-2} B_{jj}}{\min_j s_j^{-2} B_{jj}} \leq \kappa(S^{-1}BS^{-1}) \leq (d + 2) \frac{\max_j s_j^{-2} B_{jj}}{\min_j s_j^{-2} B_{jj}}$$

and, for the Jacobi preconditioning  $s_j^2 = B_{jj}$ , we have arrived at the following theorem by Wathen [19] who studies the effects of the diagonal scaling on the condition number of the Galerkin mass matrix.

**Theorem 3.2** ([19, Table 1]). *The condition number of the Jacobi preconditioned Galerkin mass matrix with a simplicial mesh has a mesh-independent bound,*

$$\kappa(S^{-1}BS^{-1}) \leq d + 2.$$

Theorems 3.1 and 3.2 show that *the mesh volume-nonuniformity has a significant effect* on the condition number of the mass matrix and that *this effect is completely eliminated* by the Jacobi preconditioning.

As we will see later in Section 5, diagonal scaling plays a similar role in reducing the effects of mesh nonuniformity on the condition number of the stiffness matrix.

---

<sup>1</sup>For example, meshes in Figure 2 have  $|K_{\max}|/|K_{\min}| \rightarrow \infty$  but  $|\omega_{\max}|/|\omega_{\min}| = \mathcal{O}(1)$ .

4. LARGEST EIGENVALUE OF THE STIFFNESS MATRIX

The following lemma is valid for any dimension.

**Lemma 4.1** (Largest eigenvalue). *The largest eigenvalue of the stiffness matrix  $A = (A_{ij})$  for the linear finite element approximation of BVP (1) is bounded by*

$$(12) \quad \max_j A_{jj} \leq \lambda_{\max}(A) \leq (d + 1) \max_j A_{jj}.$$

*The largest eigenvalue of the diagonally (Jacobi) preconditioned stiffness matrix  $S^{-1}AS^{-1}$  has a mesh-independent bound,*

$$(13) \quad 1 \leq \lambda_{\max}(S^{-1}AS^{-1}) \leq d + 1.$$

*Proof.* First, recall that for any symmetric positive semidefinite matrix  $M$ ,

$$\mathbf{v}^T M \mathbf{w} \leq \frac{1}{2} (\mathbf{v}^T M \mathbf{v} + \mathbf{w}^T M \mathbf{w}), \quad \forall \mathbf{v}, \mathbf{w} \in \mathbb{R}^{d+1}.$$

Then, using the local indices on  $K$ , the definition of  $A_{jj}$  from (6) and rearranging the sum according to the vertices, we have

$$\begin{aligned} \mathbf{u}^T A \mathbf{u} &= \int_{\Omega} \nabla u_h \cdot \mathbb{D} \nabla u_h \, d\mathbf{x} \\ &= \sum_{K \in \mathcal{T}_h} |K| \sum_{i_K, j_K=1}^{d+1} (u_{i_K} \nabla \phi_{j_K}) \cdot \mathbb{D}_K (u_{j_K} \nabla \phi_{j_K}) \\ &\leq (d + 1) \sum_{K \in \mathcal{T}_h} |K| \sum_{j_K=1}^{d+1} (u_{j_K} \nabla \phi_{j_K}) \cdot \mathbb{D}_K (u_{j_K} \nabla \phi_{j_K}) \\ &= (d + 1) \sum_j u_j^2 \sum_{K \in \omega_j} |K| \nabla \phi_j \cdot \mathbb{D}_K \nabla \phi_j \\ &= (d + 1) \sum_j u_j^2 A_{jj} \\ &\leq (d + 1) \|\mathbf{u}\|_2^2 \max_j A_{jj}. \end{aligned}$$

On the other hand, using the canonical basis vectors  $\mathbf{e}_j$  we have

$$\lambda_{\max}(A) \geq \mathbf{e}_j^T A \mathbf{e}_j = A_{jj}, \quad j = 1, \dots, N_{vi},$$

and altogether we get (12).

Using the same procedure for a diagonal scaling  $S = (s_j)$  we obtain

$$\max_j (s_j^{-2} A_{jj}) \leq \lambda_{\max}(S^{-1}AS^{-1}) \leq (d + 1) \max_j (s_j^{-2} A_{jj}).$$

For the Jacobi preconditioning we have  $s_j^2 = A_{jj}$ , which gives estimate (13). □

*Remark 4.2.* Bound (13) can also be obtained by using the unassembled form of  $A$  as shown in [19, Sect. 3]. However, the analysis employed in [19] cannot provide a lower bound on  $\lambda_{\min}(S^{-1}AS^{-1})$  other than the trivial one,  $\lambda_{\min}(S^{-1}AS^{-1}) \geq 0$ .

Although Lemma 4.1 gives a very tight bound on  $\lambda_{\max}(A)$ , it does not provide any explanation on how the mesh or the diffusion matrix affect the conditioning. We now derive a bound on  $\lambda_{\max}(A)$  in terms of mesh quantities and the diffusion matrix.

Let  $F_K: \hat{K} \rightarrow K$  be the affine mapping from the reference element  $\hat{K}$  to the mesh element  $K$ ,  $F'_K$  the Jacobian matrix of  $F_K$ ,  $j_K$  the local index of  $\phi_j$  on  $K$  and  $\hat{\phi}_{j_K} = F_K \circ \phi_{j_K}$  the corresponding basis function on  $\hat{K}$ .

Using the chain rule, we have

$$\begin{aligned}
 A_{jj} &= \sum_{K \in \omega_j} |K| \nabla \phi_j \cdot \mathbb{D}_K \nabla \phi_j \\
 &= \sum_{K \in \omega_j} |K| \left( (F'_K)^{-T} \hat{\nabla} \hat{\phi}_{j_K} \right) \cdot \mathbb{D}_K \left( (F'_K)^{-T} \hat{\nabla} \hat{\phi}_{j_K} \right) \\
 &\leq \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \|\hat{\nabla} \hat{\phi}_{j_K}\|_2^2 \\
 (14) \quad &\leq C_{\hat{\phi}} \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2,
 \end{aligned}$$

where  $C_{\hat{\phi}} = \max_{i_K=1, \dots, d+1} \|\hat{\nabla} \hat{\phi}_{i_K}\|_2^2$ . Combining this with Lemma 4.1 yields

$$(15) \quad \lambda_{\max}(A) \leq (d+1)C_{\hat{\phi}} \max_j \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2.$$

*Remark 4.3.* If we denote the maximum number of elements meeting at a mesh point by  $p_{\max}$ , then bound (15) implies

$$\lambda_{\max}(A) \leq p_{\max}(d+1)C_{\hat{\phi}} \max_K \left( |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \right),$$

which is comparable to the estimates mostly found in the literature (e.g., [5, 7, 18]). Note that both this bound and (15) are less tight than the bound (12).

*Remark 4.4.* Bound (15) implies that the scaling

$$\tilde{s}_j^2 = \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2, \quad j = 1, \dots, N_{vi}$$

will also lead to bounds similar to (13) and those in Section 5. In general,  $\tilde{s}_j^2 \geq A_{jj}$  (cf. (14)), although  $\tilde{s}_j^2 = A_{jj}$  in 1D or for a mesh that is uniform in the metric specified by  $\mathbb{D}^{-1}$  (cf. Section 4.3).

**4.1. Geometric interpretation ( $\mathbb{D} = I$ ).** For the simplest case of  $\mathbb{D} = I$ , bound (15) has a rather simple interpretation. The quantity  $\|(F'_K)^{-1}\|_2$  can be bounded by the reciprocal of the in-diameter  $h_{\min,K}$  of  $K$  [15, Lemma 5.1.2]. If we denote the average aspect of  $K$  by  $\bar{h}_K$  (i.e.,  $\bar{h}_K = |K|^{\frac{1}{d}}$ ), then we can rewrite (15) as

$$\lambda_{\max}(A) \leq C_{\hat{\phi}} \max_j \sum_{K \in \omega_j} |K| \left( \frac{1}{h_{\min,K}} \right)^2 = C_{\hat{\phi}} \max_j \sum_{K \in \omega_j} |K|^{\frac{d-2}{d}} \left( \frac{\bar{h}_K}{h_{\min,K}} \right)^2.$$

The ratio  $\bar{h}_K/h_{\min,K}$  is a measure of the *aspect ratio* of  $K$ . Thus, for the case of  $\mathbb{D} = I$ , the largest eigenvalue of  $A$  is bounded by the maximum volume-weighted element aspect ratio of the mesh. This is consistent with the observation by Shewchuk in [18] where a detailed discussion on the relation between the largest eigenvalue of the stiffness matrix and the element aspect ratio is available for the case of  $\mathbb{D} = I$  in  $d = 2$  and  $d = 3$  dimensions.



For the general case  $\mathbb{D} \neq I$ , on the other hand, it is more convenient to interpret bound (15) in terms of the *mesh quality measures* introduced in [11]. We now proceed with this.

**4.2. Mesh quality measures.** The first measure is the *alignment quality measure*, which can be simply viewed as an equivalent to the aspect ratio of  $K$  in the metric specified by  $\mathbb{D}_K^{-1}$ . It is defined as

$$Q_{ali, \mathbb{D}^{-1}}(K) = \left( \frac{\frac{1}{d} \operatorname{tr}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T})}{\det((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T})^{\frac{1}{d}}} \right)^{\frac{d}{2(d-1)}}$$

and measures how closely the principal directions of the circumscribed ellipsoid of  $K$  are aligned with the eigenvectors of  $\mathbb{D}_K$  and the semi-lengths of the principal axes are proportional to the eigenvalues [15]. Notice that

$$1 \leq Q_{ali, \mathbb{D}^{-1}}(K) < \infty.$$

In particular,  $Q_{ali, \mathbb{D}^{-1}}(K) = 1$  implies that  $K$  is equilateral in the metric  $\mathbb{D}_K^{-1}$ .

The second measure is the *equidistribution quality measure* defined as the ratio of the average element volume to the volume of  $K$ , both measured in the metric specified by  $\mathbb{D}_K^{-1}$ ,

$$(16) \quad Q_{eq, \mathbb{D}^{-1}}(K) = \frac{\frac{1}{N} \sigma_h}{|K|_{\mathbb{D}_K^{-1}}},$$

where  $|K|_{\mathbb{D}_K^{-1}} = |K| \det(\mathbb{D}_K)^{-\frac{1}{2}}$  is the volume of  $K$  with respect to  $\mathbb{D}_K^{-1}$  and

$$(17) \quad \sigma_h = \sum_{K \in \mathcal{T}_h} |K|_{\mathbb{D}_K^{-1}}.$$

The equidistribution quality measure satisfies

$$0 < Q_{eq, \mathbb{D}^{-1}}(K) < \infty \quad \text{and} \quad \frac{1}{N} \sum_{K \in \mathcal{T}_h} Q_{eq, \mathbb{D}^{-1}}^{-1}(K) = 1.$$

Notice that

$$\sigma_h = \sum_{K \in \mathcal{T}_h} |K| \det(\mathbb{D}_K)^{-\frac{1}{2}} \rightarrow \int_{\Omega} \det(\mathbb{D}(\mathbf{x}))^{-\frac{1}{2}} d\mathbf{x} = |\Omega|_{\mathbb{D}^{-1}}$$

as the mesh is being refined. As a consequence,  $\sigma_h$  can be considered as a constant.

**4.3. Geometric interpretation (general case).** Using the quality measures we can rewrite the key factor  $\|(F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}\|_2$  as

$$\begin{aligned} \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 &\leq \operatorname{tr}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}) \\ &= d Q_{ali, \mathbb{D}^{-1}}^{\frac{2(d-1)}{d}}(K) \left( |K| \det(\mathbb{D}_K)^{-\frac{1}{2}} \right)^{-\frac{2}{d}} \\ &= d \left( \frac{N}{\sigma_h} \right)^{\frac{2}{d}} \left[ Q_{ali, \mathbb{D}^{-1}}^{d-1}(K) Q_{eq, \mathbb{D}^{-1}}(K) \right]^{\frac{2}{d}} \end{aligned}$$

and, therefore,

$$(18) \quad \lambda_{\max}(A) \leq C \left( \frac{N}{\sigma_h} \right)^{\frac{2}{d}} \max_j \sum_{K \in \omega_j} |K| \left[ Q_{ali, \mathbb{D}^{-1}}^{d-1}(K) Q_{eq, \mathbb{D}^{-1}}(K) \right]^{\frac{2}{d}}.$$

Thus,  $\lambda_{\max}(A)$  is bounded by the maximum volume-weighted, combined alignment and equidistribution measure of the mesh in the metric  $\mathbb{D}_K^{-1}$ .

When a mesh is adapted to the coefficients of the BVP, i.e., it is uniform in the metric  $\mathbb{D}^{-1}$ , it will have the properties

$$(19) \quad Q_{ali, \mathbb{D}^{-1}}(K) = 1, \quad Q_{eq, \mathbb{D}^{-1}}(K) = 1, \quad \forall K \in \mathcal{T}_h$$

and

$$(20) \quad \left(\frac{N}{\sigma_h}\right)^{\frac{2}{d}} \leq \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \leq d \left(\frac{N}{\sigma_h}\right)^{\frac{2}{d}}.$$

Moreover, bound (18) will reduce to

$$\lambda_{\max}(A) \leq CN^{\frac{2}{d}} |\omega_{\max}|.$$

5. SMALLEST EIGENVALUE AND CONDITION NUMBER OF THE STIFFNESS MATRIX

The approach employed in this section was originally developed by Bank and Scott [2] for isotropic meshes. We generalize it here to arbitrary anisotropic meshes.

Hereafter, we will use  $C$  as a generic constant which can have different values at different appearances but is independent of the mesh, the number of mesh elements, and the solution of the BVP.

We start with bounds on  $\lambda_{\min}(A)$ .

**Lemma 5.1** (Smallest eigenvalue). *The smallest eigenvalue of the stiffness matrix for the linear finite element approximation of BVP (1) is bounded from below by*

$$(21) \quad \lambda_{\min}(A) \geq Cd_{\min} N^{-1} \begin{cases} 1, & \text{for } d = 1, \\ \left(1 + \ln \frac{|\bar{K}|}{|K_{\min}|}\right)^{-1}, & \text{for } d = 2, \\ \left(\frac{1}{N} \sum_{K \in \mathcal{T}_h} \left(\frac{|\bar{K}|}{|K|}\right)^{\frac{d-2}{2}}\right)^{-\frac{2}{d}}, & \text{for } d \geq 3, \end{cases}$$

where  $|\bar{K}| = \frac{1}{N} |\Omega|$  denotes the average element size.

The smallest eigenvalue of the diagonally (Jacobi) preconditioned stiffness matrix is bounded from below by

$$(22) \quad \lambda_{\min}(S^{-1}AS^{-1}) \geq CN^{-2} \left(\frac{1}{Nd_{\min}} \sum_{K \in \mathcal{T}_h} \mathbb{D}_K \frac{|\bar{K}|}{|K|}\right)^{-1}, \quad \text{for } d = 1$$

and

$$(23) \quad \lambda_{\min}(S^{-1}AS^{-1}) \geq CN^{-\frac{2}{d}} \left(\frac{1}{Nd_{\min}^{\frac{d}{2}}} \sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2^{\frac{d}{2}}\right)^{-\frac{2}{d}} \\ \times \begin{cases} \left(1 + \left| \ln \frac{\max_{K \in \mathcal{T}_h} \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2}{\sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2} \right| \right)^{-1}, & \text{for } d = 2, \\ 1, & \text{for } d \geq 3. \end{cases}$$

*Proof.* Since Sobolev’s inequality is different for  $d = 1$ ,  $d = 2$  and  $d \geq 3$  dimensions [8, Theorem 7.10], we treat these cases separately.

Case  $d = 1$ . Let  $C_S$  be the constant associated with Sobolev’s inequality. Using the inequality (2), Sobolev’s inequality, and the equivalence of the vector norms,

$$\begin{aligned} \mathbf{u}^T \mathbf{A} \mathbf{u} &\geq d_{\min} |u_h|_{H^1(\Omega)}^2 \\ &\geq d_{\min} C_S |\Omega|^{-1} \sup_{\Omega} |u_h|^2 \\ &= d_{\min} C_S |\Omega|^{-1} \max_j u_j^2 \\ &\geq d_{\min} C_S |\Omega|^{-1} N^{-1} \|\mathbf{u}\|_2^2. \end{aligned}$$

Therefore,  $\lambda_{\min}(A) \geq C d_{\min} N^{-1}$ .

With scaling,

$$(24) \quad \mathbf{u}^T S^{-1} A S^{-1} \mathbf{u} \geq C d_{\min} \max_j s_j^{-2} u_j^2 \geq C d_{\min} \frac{\sum_j s_j^2 s_j^{-2} u_j^2}{\sum_j s_j^2} = C d_{\min} \frac{\|\mathbf{u}\|_2^2}{\sum_j s_j^2}.$$

In 1D,  $\nabla \phi_{jK}|_K = \pm |K|^{-1}$  and, therefore,

$$s_j^2 = A_{jj} = \sum_{K \in \omega_j} |K| \nabla \phi_j \cdot \mathbb{D}_K \nabla \phi_j = \sum_{K \in \omega_j} \frac{\mathbb{D}_K}{|K|}.$$

Using this in (24) gives (22).

Case  $d = 2$ . Consider a set of not-all-zero nonnegative numbers  $\{\alpha_K, K \in \mathcal{T}_h\}$  (to be determined) and a finite number  $q > 2$ . Let  $C_P$ ,  $C_S$ , and  $C_{\hat{K}}$  be the constants associated with Poincaré’s inequality, Sobolev’s inequality, and the norm equivalence on  $\hat{K}$ , respectively. Using (2), Poincaré’s, Sobolev’s and Hölder’s inequalities and the norm equivalence for  $\hat{u}_h$ , we have

$$\begin{aligned} \mathbf{u}^T \mathbf{A} \mathbf{u} &= \int_{\Omega} \nabla u_h \cdot \mathbb{D} \nabla u_h \, d\mathbf{x} \geq d_{\min} |u_h|_{H^1(\Omega)}^2 \\ &\geq \frac{d_{\min} C_P}{1 + C_P} \|u_h\|_{H^1(\Omega)}^2 \\ &\geq \frac{d_{\min} C_P C_S}{1 + C_P} \frac{1}{q} \|u_h\|_{L^q(\Omega)}^2 \\ &= \frac{d_{\min} C_P C_S}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \|u_h\|_{L^q(K)}^q \right)^{\frac{2}{q}} \\ &= \frac{d_{\min} C_P C_S}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{\frac{q-2}{q}} \left( \sum_{K \in \mathcal{T}_h} \|u_h\|_{L^q(K)}^q \right)^{\frac{2}{q}} \\ &\geq \frac{d_{\min} C_P C_S}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_{K \in \mathcal{T}_h} \alpha_K \|u_h\|_{L^q(K)}^2 \\ &= \frac{d_{\min} C_P C_S}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_{K \in \mathcal{T}_h} \alpha_K |K|^{\frac{2}{q}} \|\hat{u}_h\|_{L^q(\hat{K})}^2 \end{aligned}$$

$$\begin{aligned} &\geq \frac{d_{\min} C_P C_S C_{\hat{K}}}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_{K \in \mathcal{T}_h} \alpha_K |K|^{\frac{2}{q}} \|\mathbf{u}_K\|_2^2 \\ &= \frac{d_{\min} C_P C_S C_{\hat{K}}}{1 + C_P} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_j u_j^2 \sum_{K \in \omega_j} \alpha_K |K|^{\frac{2}{q}}. \end{aligned}$$

The choice  $\alpha_K = |K|^{-\frac{2}{q}}$  yields

$$\mathbf{u}^T A \mathbf{u} \geq C d_{\min} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} |K|^{-\frac{2}{q-2}} \right)^{-\frac{q-2}{q}} \sum_j u_j^2$$

and, therefore,

$$\begin{aligned} \lambda_{\min}(A) &\geq C d_{\min} q^{-1} \left( \sum_{K \in \mathcal{T}_h} |K|^{-\frac{2}{q-2}} \right)^{-\frac{q-2}{q}} \\ &\geq C d_{\min} q^{-1} \left( N |K_{\min}|^{-\frac{2}{q-2}} \right)^{-\frac{q-2}{q}} \\ (25) \quad &= C d_{\min} N^{-1} \left[ q^{-1} (N |K_{\min}|)^{\frac{2}{q}} \right]. \end{aligned}$$

The largest lower bound on (25) is obtained for  $q = \max \{2, |\ln(N |K_{\min}|)|\}$  with

$$q^{-1} (N |K_{\min}|)^{\frac{2}{q}} \geq \frac{C}{1 + |\ln(N |K_{\min}|)|}.$$

The choice  $q = 2$  is viewed as the limiting case as  $q \rightarrow 2^+$ . Estimate (21) follows from this, (25) and the definition of the average element size.

With scaling, we have

$$\mathbf{u}^T S^{-1} A S^{-1} \mathbf{u} \geq C d_{\min} \frac{1}{q} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_j u_j^2 s_j^{-2} \sum_{K \in \omega_j} \alpha_K |K|^{\frac{2}{q}}.$$

For the Jacobi preconditioning  $s_j^2 = A_{jj} = \sum_{K \in \omega_j} \nabla \phi_j \cdot \mathbb{D}_K \nabla \phi_j$  we choose

$$\alpha_K = |K|^{\frac{q-2}{q}} \sum_{i_K=1}^{d+1} \nabla \phi_{i_K} \cdot \mathbb{D}_K \nabla \phi_{i_K} = |K|^{\frac{q-2}{q}} \sum_{i_K=1}^{d+1} \hat{\nabla} \hat{\phi}_{i_K} \cdot (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \hat{\nabla} \hat{\phi}_{i_K},$$

which gives

$$s_j^{-2} \sum_{K \in \omega_j} \alpha_K |K|^{\frac{2}{q}} \geq 1$$

and

$$\alpha_K \leq (d + 1) C_{\hat{\phi}} |K|^{\frac{q-2}{q}} \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2,$$

where  $C_{\hat{\phi}} = \max_{i_K=1, \dots, d+1} \|\hat{\nabla} \hat{\phi}_{i_K}\|^2$ . With these and choosing the value for the index  $q$  in a similar manner as for the case without scaling we obtain (23).

*Case  $d \geq 3$ .* This case is very similar to case  $d = 2$ . Again, from (2), Poincaré’s, Sobolev’s and Hölder’s inequalities and the norm equivalence for  $\hat{u}_h$ , we have

$$\begin{aligned}
 \mathbf{u}^T A \mathbf{u} &= \int_{\Omega} \nabla u_h \cdot \mathbb{D} \nabla u_h \, d\mathbf{x} \geq d_{\min} \|u_h\|_{H^1(\Omega)}^2 \\
 &\geq \frac{d_{\min} C_P}{1 + C_P} \|u_h\|_{H^1(\Omega)}^2 \\
 &\geq \frac{d_{\min} C_P C_S}{1 + C_P} \|u_h\|_{L^{\frac{2d}{d-2}}(\Omega)}^2 \\
 &= \frac{d_{\min} C_P C_S}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \|u_h\|_{L^{\frac{2d}{d-2}}(K)} \right)^{\frac{d-2}{d}} \\
 &= \frac{d_{\min} C_P C_S}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{\frac{2}{d}} \left( \sum_{K \in \mathcal{T}_h} \|u_h\|_{L^{\frac{2d}{d-2}}(K)} \right)^{\frac{d-2}{d}} \\
 &\geq \frac{d_{\min} C_P C_S}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \sum_{K \in \mathcal{T}_h} \alpha_K \|u_h\|_{L^{\frac{2d}{d-2}}(K)}^2 \\
 &= \frac{d_{\min} C_P C_S}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \sum_{K \in \mathcal{T}_h} \alpha_K |K|^{\frac{d-2}{d}} \|\hat{u}_h\|_{L^{\frac{2d}{d-2}}(\hat{K})}^2 \\
 &\geq \frac{d_{\min} C_P C_S C_{\hat{K}}}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \sum_{K \in \mathcal{T}_h} \alpha_K |K|^{\frac{d-2}{d}} \|\mathbf{u}_K\|_2^2 \\
 &= \frac{d_{\min} C_P C_S C_{\hat{K}}}{1 + C_P} \left( \sum_{K \in \mathcal{T}_h} \alpha_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \sum_j u_j^2 \sum_{K \in \omega_j} \alpha_K |K|^{\frac{d-2}{d}}.
 \end{aligned}$$

The choice  $\alpha_K = |K|^{-\frac{d-2}{d}}$  gives

$$\mathbf{u}^T A \mathbf{u} \geq C d_{\min} \left( \sum_{K \in \mathcal{T}_h} |K|^{\frac{2-d}{2}} \right)^{-\frac{2}{d}} \sum_j u_j^2.$$

Estimate (21) follows from this and the definition of the average element size.

The bound for the scaled stiffness matrix is obtained by choosing

$$\alpha_K = |K|^{\frac{2}{d}} \sum_{i_K=1}^{d+1} \hat{\nabla} \hat{\phi}_{i_K} \cdot (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \hat{\nabla} \hat{\phi}_{i_K}. \quad \square$$

Combining Lemma 4.1, estimate (15) and Lemma 5.1 we obtain upper bounds on the condition number of the stiffness matrix and the scaled stiffness matrix.

**Theorem 5.2** (Condition number of the stiffness matrix). *The condition number of the stiffness matrix for the linear finite element approximation of BVP (1) is bounded by*

$$(26) \quad \kappa(A) \leq C N^2 \frac{1}{d_{\min}} \max_j \sum_{K \in \omega_j} \mathbb{D}_K \frac{|\bar{K}|}{|K|}, \quad \text{for } d = 1$$

and

$$(27) \quad \kappa(A) \leq CN^{\frac{2}{d}} \left( \frac{N^{1-\frac{2}{d}}}{d_{\min}} \max_j \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \right) \times \begin{cases} 1 + \ln \frac{|\bar{K}|}{|K_{\min}|}, & \text{for } d = 2, \\ \left( \frac{1}{N} \sum_{K \in \mathcal{T}_h} \left( \frac{|\bar{K}|}{|K|} \right)^{\frac{d-2}{2}} \right)^{\frac{2}{d}}, & \text{for } d \geq 3. \end{cases}$$

The condition number of the diagonally (Jacobi) preconditioned stiffness matrix is bounded by

$$(28) \quad \kappa(S^{-1}AS^{-1}) \leq CN^2 \frac{1}{Nd_{\min}} \sum_{K \in \mathcal{T}_h} \mathbb{D}_K \frac{|\bar{K}|}{|K|}, \quad \text{for } d = 1$$

and

$$(29) \quad \kappa(S^{-1}AS^{-1}) \leq CN^{\frac{2}{d}} \left( \frac{1}{Nd_{\min}^{\frac{d}{2}}} \sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2^{\frac{d}{2}} \right) \times \begin{cases} 1 + \left| \ln \frac{\max_{K \in \mathcal{T}_h} \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2}{\sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2} \right|, & \text{for } d = 2, \\ 1, & \text{for } d \geq 3. \end{cases}$$

We now study the geometric interpretation of the bounds.

**5.1. Geometric interpretation (without scaling).** Bounds (26) and (27) contain three factors, a base bound  $CN^{\frac{2}{d}}$ , a factor reflecting the effects of the mesh nonuniformity measured in the metric  $\mathbb{D}^{-1}$  (*mesh  $\mathbb{D}$ -nonuniformity*), and, if  $d \geq 2$ , a factor reflecting the effects of the mesh nonuniformity in volume measured in the Euclidean metric (*volume-nonuniformity*).

The first factor  $N^{\frac{2}{d}}$  corresponds to the condition number of the stiffness matrix for the Laplacian operator on a uniform mesh (cf. Special Case 5.1 below).

The second factor

$$\frac{N^{1-\frac{2}{d}}}{d_{\min}} \max_j \sum_{K \in \omega_j} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2$$

reflects the effects of the mesh  $\mathbb{D}$ -nonuniformity and can be understood as a volume-weighted, combined alignment and equidistribution quality measure of the mesh with respect to  $\mathbb{D}^{-1}$  (cf. Section 4.3).

The third factor in (27) is

$$\begin{cases} 1 + \ln \frac{|\bar{K}|}{|K_{\min}|}, & \text{for } d = 2, \\ \left( \frac{1}{N} \sum_{K \in \mathcal{T}_h} \left( \frac{|\bar{K}|}{|K|} \right)^{\frac{d}{2}-1} \right)^{\frac{2}{d}}, & \text{for } d \geq 3. \end{cases}$$

It measures the effects of the mesh volume-nonuniformity (measured in the Euclidean metric) on the condition number. Notice that there is no effect in 1D and in 2D it is minimal. In  $d \geq 3$  dimensions the factor is proportional to the

average of  $|K|^{-1+\frac{2}{d}}$  over all elements. This is a significant improvement in comparison with previously available estimates which are proportional to  $|K_{\min}|^{-1}$  [6] or  $|K_{\min}|^{-1+\frac{2}{d}}$  [7].

**5.2. Geometric interpretation (with scaling).** Bounds (28) and (29) for the scaled stiffness matrix have the same base bound as without scaling. Hence, diagonal scaling has no effect on the condition number when the mesh is uniform and  $\mathbb{D} = I$ .

Unlike (27), bounds (28) and (29) do not have the third factor which involves only the element volume (in comparison to the second factor which couples  $(F'_K)^{-1}$  with  $\mathbb{D}^{-1}$ ). In this sense, a properly chosen diagonal scaling can eliminate the effects of the mesh volume-nonuniformity on the condition number. Moreover, scaling can also significantly reduce the effects of the mesh  $\mathbb{D}$ -nonuniformity. Indeed, the factors in (28) and (29) that couple  $(F'_K)^{-1}$  with  $\mathbb{D}^{-1}$  are asymptotically the  $L^{\max\{1, \frac{d}{2}\}}(\Omega)$  norm of  $\|(F'_K)^{-1}\mathbb{D}_K(F'_K)^{-T}\|_2$  whereas the corresponding factors in (26) and (27) are basically the maximum norm.

Furthermore, the  $\mathbb{D}$ -related factor in (29) for  $d \geq 2$  can be rewritten in terms of the alignment quality measure  $Q_{ali, \mathbb{D}^{-1}}$  from Section 4.2 as

$$(30) \quad \sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2^{\frac{d}{2}} \leq d^{\frac{d}{2}} \sum_{K \in \mathcal{T}_h} Q_{ali, \mathbb{D}^{-1}}^{d-1}(K) \det(\mathbb{D}_K)^{\frac{1}{2}}.$$

Thus, the dependence of this  $\mathbb{D}$ -related factor on the element volume is also mild: both  $Q_{ali, \mathbb{D}^{-1}}(K)$  and  $\mathbb{D}_K$  (the average of  $\mathbb{D}$  over  $K$ ) are invariant under the scaling transformation of  $K$ .

The following special cases are instructional to understand the interplay of the factors for different types of meshes.

**Special Case 5.1** (Uniform meshes). For a uniform mesh and  $\mathbb{D} = I$ , bounds (26)–(29) yield

$$\kappa(A) \leq CN^{\frac{2}{d}} \quad \text{and} \quad \kappa(S^{-1}AS^{-1}) \leq CN^{\frac{2}{d}},$$

which is the base bound. Hence, the diagonal scaling has no effect on the condition number when the mesh is uniform and  $\mathbb{D} = I$ .

**Special Case 5.2** (Isotropic meshes,  $\mathbb{D} = I$ ,  $d \geq 2$ ). For an isotropic mesh and  $\mathbb{D} = I$ ,

$$|K| \sim h_K^d \quad \text{and} \quad \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \sim h_K^{-2}.$$

Therefore,

$$\frac{1}{N} \sum_{K \in \mathcal{T}_h} |K| \left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2^{\frac{d}{2}} \lesssim \frac{1}{N} \sum_{K \in \mathcal{T}_h} h_K^d h_K^{-d} = 1$$

and bound (29) reduces to

$$(31) \quad \kappa(S^{-1}AS^{-1}) \leq CN^{\frac{2}{d}} \begin{cases} 1 + \ln \frac{|\bar{K}|}{|K_{\min}|}, & \text{for } d = 2, \\ 1, & \text{for } d \geq 3, \end{cases}$$

which is precisely the result of Bank and Scott [2, Theorems 4.2 and 5.2]. In this case, the diagonal scaling becomes

$$s_j = (A_{jj})^{\frac{1}{2}} = \left( \sum_{K \in \omega_j} |K| \nabla \phi_j \cdot \nabla \phi_j \right)^{\frac{1}{2}} \sim \left( \sum_{K \in \omega_j} h_K^{d-2} \right)^{\frac{1}{2}} \sim h_j^{\frac{d-2}{2}},$$

where  $h_j$  denotes the average length of the elements around the  $j^{\text{th}}$  vertex. This scaling is equivalent to the change of basis functions

$$\phi_j \rightarrow h_j^{\frac{2-d}{2}} \phi_j,$$

which is used in [2, Example 2.1].

**Special Case 5.3** (Uniform meshes with respect to  $\mathbb{D}^{-1}$ ). For a mesh that is uniform with respect to  $\mathbb{D}^{-1}$ , i.e., *coefficient adaptive*, we have properties (19) and (20). Bounds (26)–(29) reduce to

$$\kappa(A) \leq \frac{C(N |\omega_{\max}|)}{d_{\min}} \left( \frac{N}{\sigma_h} \right)^{\frac{2}{d}} \begin{cases} 1, & \text{for } d = 1, \\ 1 + \ln \frac{|\bar{K}|}{|K_{\min}|}, & \text{for } d = 2, \\ \left( \frac{1}{N} \sum_{K \in \mathcal{T}_h} \left( \frac{|\bar{K}|}{|K|} \right)^{\frac{d-1}{2}} \right)^{\frac{2}{d}}, & \text{for } d \geq 3, \end{cases}$$

$$\kappa(S^{-1}AS^{-1}) \leq \frac{C}{d_{\min}} \left( \frac{N}{\sigma_h} \right)^{\frac{2}{d}}, \quad \text{for } d \geq 1,$$

where  $\sigma_h$  is defined in (17) and corresponds to the volume of the domain in the metric specified by  $\mathbb{D}^{-1}$ . Thus, the condition number of the scaled stiffness matrix for a coefficient adaptive mesh has the optimal order of  $\mathcal{O}(N^{\frac{2}{d}})$ .

**Special Case 5.4** (Aligned meshes,  $d \geq 2$ ). For meshes aligned with the diffusion matrix but not necessarily fully coefficient adaptive (i.e., isotropic but not uniform with respect to  $\mathbb{D}^{-1}$ ) we have

$$Q_{ali, \mathbb{D}^{-1}}(K) = 1 \quad \text{but} \quad Q_{eq, \mathbb{D}^{-1}}(K) \neq 1.$$

From (30), bound (29) becomes

$$\kappa(S^{-1}AS^{-1}) \leq C \frac{N^{\frac{2}{d}}}{d_{\min}} \left( \frac{1}{N} \sum_{K \in \mathcal{T}_h} \det(\mathbb{D}_K)^{\frac{1}{2}} \right)^{\frac{2}{d}} \begin{cases} 1 + \ln \frac{|\bar{K}|}{|K_{\min}|}, & \text{for } d = 2, \\ 1, & \text{for } d \geq 3. \end{cases}$$

Aside from the term depending on  $\det(\mathbb{D})$ , this bound is equivalent to (31). Hence, the diagonal scaling almost eliminates the effects of the mesh on the condition number for  $\mathbb{D}$ -aligned meshes.

**Special Case 5.5** (General  $M$ -uniform meshes). Finally, let us consider general  $M$ -uniform meshes, i.e., meshes that are uniform in the metric specified by a given metric tensor  $M$  which does not necessarily correspond to  $\mathbb{D}^{-1}$ . In the context of mesh adaptation, an adaptive mesh is typically generated based on some estimate of the solution error and the associated metric tensor  $M$  is solution dependent. Thus, it is of interest to know what the impact of a given  $M$  on the conditioning of the stiffness matrix is. Recall [13] that an  $M$ -uniform mesh satisfies

$$(F'_K)^{-T} (F'_K)^{-1} = \left( \frac{N}{\sigma_{h,M}} \right)^{\frac{2}{d}} M_K,$$



where  $M_K$  is some average of  $M$  on  $K$  and  $\sigma_{h,M}$  is defined as in (17) but with  $\mathbb{D}$  replaced by  $M^{-1}$ . We have

$$\left\| (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right\|_2 \leq \left( \frac{N}{\sigma_{h,M}} \right)^{\frac{2}{d}} \|M_K \mathbb{D}_K\|_2$$

and, therefore,

$$\kappa(S^{-1}AS^{-1}) \leq \frac{C}{d_{\min}} \left( \frac{N}{\sigma_{h,M}} \right)^{\frac{2}{d}} \left( \sum_{K \in \mathcal{T}_h} |K| \|M_K \mathbb{D}_K\|_2^{\frac{d}{2}} \right)^{\frac{2}{d}}.$$

Hence, the bound on the condition number after diagonal scaling for an  $M$ -uniform mesh depends only on the volume-weighted average of  $\|M_K \mathbb{D}_K\|_2^{d/2}$  or, asymptotically, the  $L^{d/2}$  norm of  $\|M \mathbb{D}\|_2$ . For many problems such as those having boundary layers and shock waves, mesh elements are typically concentrated in a small portion of the physical domain. In that situation, we would expect that  $M$  differs significantly from  $\mathbb{D}^{-1}$  only in small regions. As a consequence, the volume-weighted average of  $\|M_K \mathbb{D}_K\|_2^{d/2}$  over the whole domain may remain small and, therefore, the condition number of the scaled stiffness matrix for anisotropic adaptive meshes does not necessarily increase as much as generally assumed.

This effect can be observed in Examples 6.2 and 6.3. Figures 3a and 4a show that the effects of anisotropic adaptation are completely neutralized by the diagonal scaling when the number of anisotropic elements is small in comparison to  $N$ .

### 6. NUMERICAL EXPERIMENTS

In this section we present numerical results for a selection of one-, two-, and three-dimensional examples to illustrate our theoretical findings.

Note that all bounds on the smallest eigenvalue contain a constant  $C$ . We obtain its value by calibrating the bound for  $\lambda_{\min}(S^{-1}AS^{-1})$  with Delaunay (Example 6.4) or uniform meshes (all other examples) through comparing the exact and estimated values. For the largest eigenvalue we use explicit bounds (12) and (13).

First, we give examples with predefined meshes to demonstrate the influence of the number and shape of mesh elements on the condition number of the stiffness matrix and to verify the improvement achieved with the diagonal scaling. For the tests, we employ the Laplace operator (i.e.,  $\mathbb{D} = I$ ) and a mesh on the unit interval, square, and cube, for 1D, 2D, and 3D, respectively.

**Example 6.1** ( $d = 1$ ,  $\mathbb{D} = I$ , Chebyshev nodes). For a simple one-dimensional example we choose a mesh given by Chebyshev nodes in the interval  $[0, 1]$ ,

$$(32) \quad x_i = \frac{1}{2} \left( 1 - \cos \frac{\pi(2i - 1)}{2(N - 1)} \right), \quad i = 1, \dots, N - 1.$$

The exact condition number of the stiffness matrix and its estimates (26) and (28) are shown in Figures 1a (without scaling) and 1b (with scaling) while those for the extreme eigenvalues and their estimates are given in Figures 1c (without scaling) and 1d (with scaling).

Figure 1a shows that the estimate (26) is much sharper than the standard estimate with  $\lambda_{\min}(A) \propto |K_{\min}|$ . The former has the same asymptotic order as the exact value as  $N$  increases, whereas the latter is too pessimistic and has a higher asymptotic order. The difference is caused by the estimate of the smallest eigenvalue

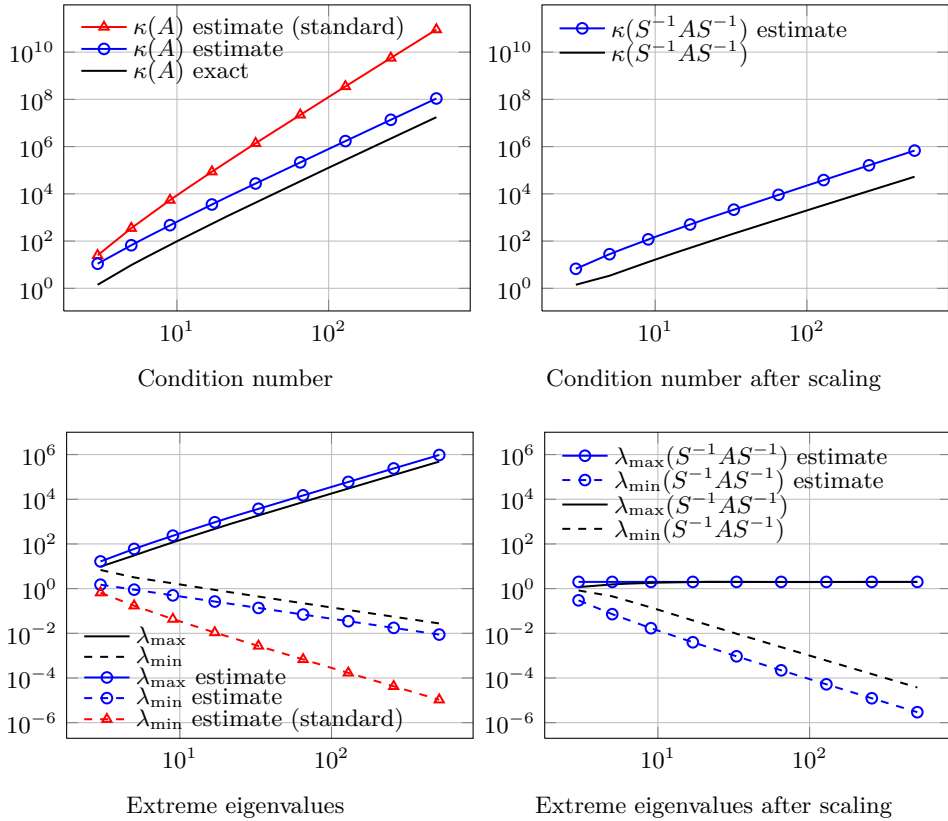


FIGURE 1. Example 6.1: Exact and estimated condition number and eigenvalues of the stiffness matrix as a function of  $N$  ( $d = 1$ )

(Figure 1c). Notice that the estimates on the largest eigenvalue are very tight, both for the scaled and the unscaled cases.

The results clearly show the benefits of diagonal scaling: the order for the condition number of the scaled stiffness matrix in Figure 1b is  $\mathcal{O}(N^2 \ln N)$ , which is almost the same as for uniform meshes, whereas that without scaling in Figure 1a is  $\mathcal{O}(N^3)$ . It can be shown analytically that the orders of the nonuniformity factors in (26) and (28) for the Chebyshev nodes defined with (32) are  $\mathcal{O}(N)$  and  $\mathcal{O}(\ln N)$  and those of the corresponding condition numbers are  $\mathcal{O}(N^3)$  and  $\mathcal{O}(N^2 \ln N)$ .

Thus, the numerical and theoretical results are consistent and the improvement by diagonal scaling from the maximum norm to the  $L^2$  norm is significant in this example.

**Example 6.2** ( $d = 2$ ,  $\mathbb{D} = I$ , anisotropic elements in a unit square). For this 2D example we use a mesh for the unit square  $[0, 1] \times [0, 1]$  with  $\mathcal{O}(N^{1/2})$  skew elements, as shown in Figure 2a. First, we fix the maximum aspect ratio at 125 : 1 and increase  $N$  to verify the dependence of the condition number on  $N$  (Figure 3a). Then, we fix  $N$  at 20,000 and change the maximum aspect ratio of the mesh elements to investigate the dependence of the conditioning on the mesh shape (Figure 3b).

Figure 3a shows the averaging effect of the diagonal scaling: the scaling significantly reduces the condition number and, when  $N$  becomes large enough, the conditioning of a scaled system is comparable to the condition number on a uniform mesh. Moreover, the estimated value of the condition number with or without scaling has the same order as the exact value as  $N$  increases.

Figure 3b provides a good numerical validation of (27), namely that the condition number of the unscaled stiffness matrix is linearly proportional to the largest aspect ratio.<sup>2</sup> With scaling, the condition number is still increasing with an increasing aspect ratio, since the average aspect ratio is also increasing (in accordance to (29)). Nevertheless, the condition number after scaling is smaller by a factor of 10.

Figure 3b also shows that our estimate of the condition number with scaling has the same (linear) order as the exact value as the maximum aspect ratio increases, whereas the bounds for the unscaled case has a slightly higher order. This indicates that the estimation can be further improved.

As for the estimates on the extreme eigenvalues, the results are mainly the same as in Example 6.1. For this reason, we omit them in 2D and 3D to save space.

**Example 6.3** ( $d = 3$ ; anisotropic elements in a unit cube). In this example, we repeat the same test setting as in Example 6.2: fixed anisotropy (25 : 25 : 1) with increasing number of elements (Figure 4a) and a fixed  $N = 29,478$  paired with the changing anisotropy of the mesh (Figure 4b). The results shown in Figure 4 are essentially the same as in 2D. Since the mesh used in this example has a larger share of skew elements ( $\mathcal{O}(N^{-1/3})$ ) than the mesh used in Example 6.2 ( $\mathcal{O}(N^{-1/2})$ ), it is reasonable to expect that the averaging effect of diagonal scaling is less effective. This can be seen in Figure 4 where the exact condition numbers with and without scaling stay closer than in Figure 3.

Figure 4 shows that the bounds on the condition number with and without scaling have the same asymptotic order as the exact values as  $N$  increases. However, they have higher orders as the maximum aspect ratio increases for a fixed  $N$ . As in the previous example, this indicates that the estimation can be further improved.

In the next example, we consider an adaptive finite element solution of an anisotropic diffusion problem with different meshes.

**Example 6.4** ( $d = 2$ , adaptive anisotropic meshes). Consider an anisotropic diffusion problem studied in [16, 17]. It takes the form of BVP (1) but with a non-homogeneous Dirichlet boundary condition. The domain and its outer and inner boundaries  $\Gamma_{\text{out}}$  and  $\Gamma_{\text{in}}$  are shown in Figure 5a. The coefficients of the BVP, the right-hand side and the boundary data are given by

$$(33) \quad \mathbb{D} = \begin{pmatrix} \cos \psi & -\sin \psi \\ \sin \psi & \cos \psi \end{pmatrix} \begin{pmatrix} 1000 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \psi & \sin \psi \\ -\sin \psi & \cos \psi \end{pmatrix}, \quad \psi = \pi \sin x \cos y,$$

$$f = 0 \text{ in } \Omega = (0, 1)^2 \setminus \left[ \frac{4}{9}, \frac{5}{9} \right]^2, \quad g = 0 \text{ on } \Gamma_{\text{out}}, \quad g = 2 \text{ on } \Gamma_{\text{in}}.$$

We employ an adaptive finite element algorithm from [14, 16] to compute the numerical solution and adaptive meshes. The algorithm utilizes the  $M$ -uniform mesh approach, i.e., meshes are generated as quasi-uniform in a given metric  $M$ . For the mesh generation we use the *bidimensional anisotropic mesh generator* [9].

<sup>2</sup>In 2D with  $\mathbb{D} = I$ , the nonuniformity term in (27) is equivalent to the aspect ratio.

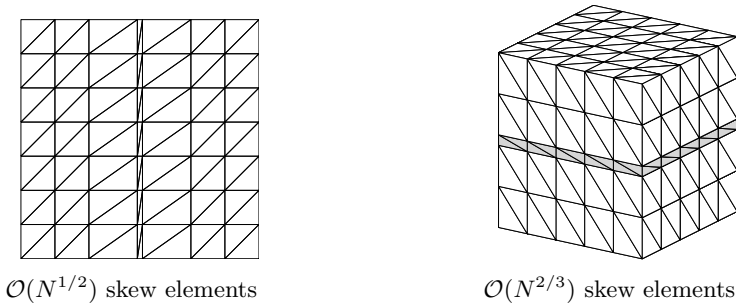


FIGURE 2. Predefined meshes for (a) Example 6.2 and (b) Example 6.3

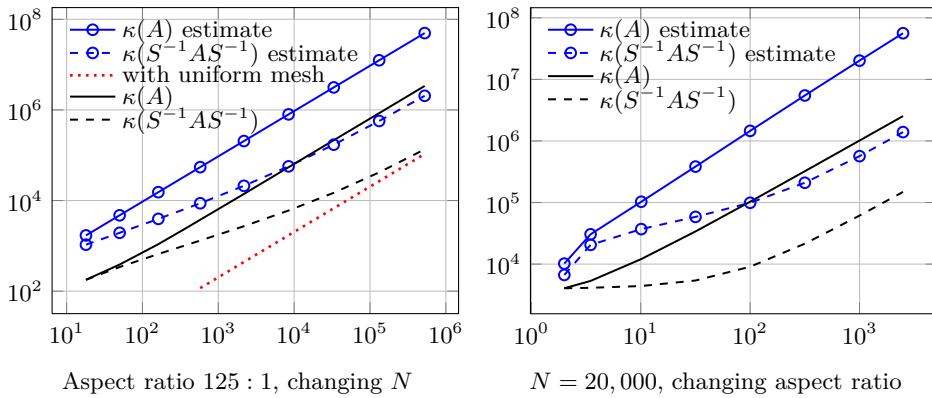


FIGURE 3. Example 6.2: Condition number before and after scaling for a predefined 2D mesh (Figure 2a) as a function of (a) the number of mesh elements and (b) the maximum element aspect ratio

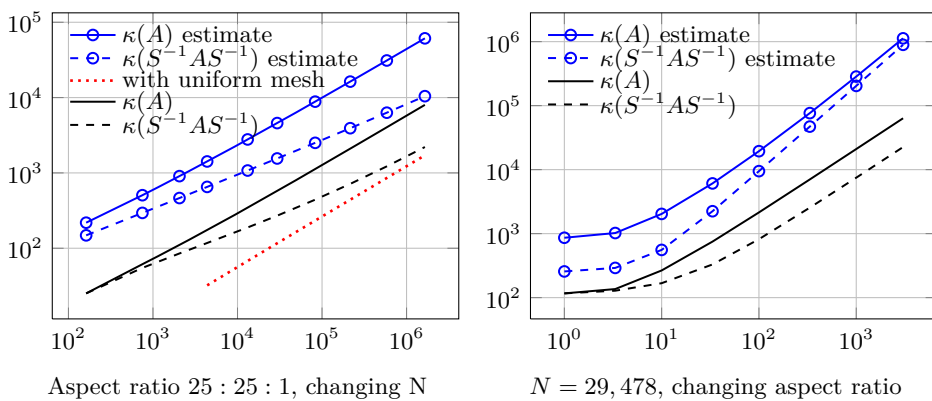


FIGURE 4. Example 6.3: Condition number before and after scaling for a predefined 3D mesh (Figure 2b) as a function of (a) the number of mesh elements and (b) the maximum element aspect ratio

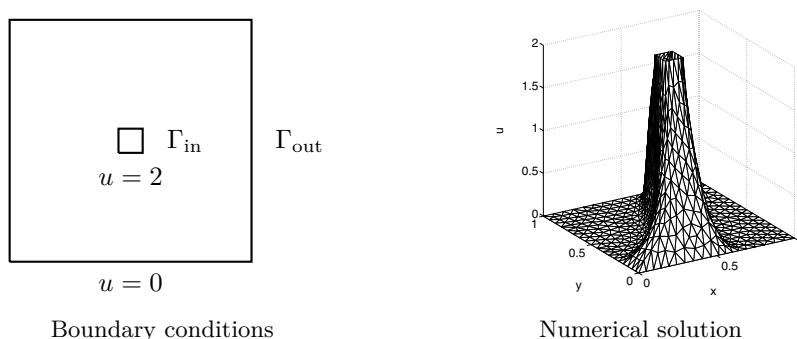


FIGURE 5. Example 6.4: Boundary conditions and a numerical solution

A Delaunay mesh—our first example (Figure 6a)—is  $M$ -uniform (or  $M$ -quasi-uniform) with respect to  $M = I$ . The second mesh (Figure 6b) is purely coefficient-adaptive and is defined as an  $M$ -uniform mesh with respect to  $\mathbb{D}^{-1}$ , i.e.,  $M = \mathbb{D}^{-1}$ . The third mesh (Figure 6c) is a purely solution-adaptive mesh where  $M = M(u_h)$  depends on the numerical solution  $u_h$  (or, more precisely, on the hierarchical basis error estimate  $e_h$ ). The fourth mesh (Figure 6d) represents a combination of adaptation to both the solution and the coefficients of the problem and the metric is defined as  $M = \theta(e_h)\mathbb{D}^{-1}$ , where  $\theta(e_h)$  is a scalar function depending on the error estimator  $e_h$ . With such choice the shape of mesh elements is determined by the diffusion matrix while the size is controlled by the estimate of the solution error.

From Figure 6 we can see that the smallest condition number among all four meshes is with the purely coefficient-adaptive mesh (Figure 6b), which is consistent with Special Case 5.3. The conditioning is better than in the case of a quasi-uniform mesh (Figure 6a), confirming the observation that, depending on the problem, a quasi-uniform mesh is not necessarily the best mesh from the conditioning point of view. For both cases, diagonal scaling does not improve the condition number significantly. This is expected since both meshes are almost volume-uniform. To explain why the mesh in Figure 6b is (almost) volume-uniform, we recall from (16) and (19) that an  $M$ -uniform mesh with respect to  $\mathbb{D}^{-1}$  satisfies

$$|K| \sim \sqrt{\det(\mathbb{D}_K)}, \quad \forall K \in \mathcal{T}_h.$$

The diffusion matrix  $\mathbb{D}$  in (33) satisfies  $\det(\mathbb{D}) = 1000$ . Thus,  $|K| = \text{const}$ .

The largest condition number is in the case of the purely solution-adaptive mesh (Figure 6c). This is because the mesh is not volume-uniform and its elements are not aligned with  $\mathbb{D}^{-1}$ . Since the mesh is far from being uniform in size, scaling will have a significant impact, as can be verified in Figure 6c: the condition number after the scaling is even smaller than the condition number with Delaunay meshes.

Conditioning with a mesh that is both coefficient- and solution-adaptive (Figure 6d) is not as good as in the case of the purely coefficient-adaptive mesh but better than in the case of the purely adaptive and Delaunay meshes.

In all four cases we observe that the developed estimates for the condition number of the stiffness matrix are reasonably tight and have the same order as the exact values as  $N$  increases for both unscaled and scaled cases.

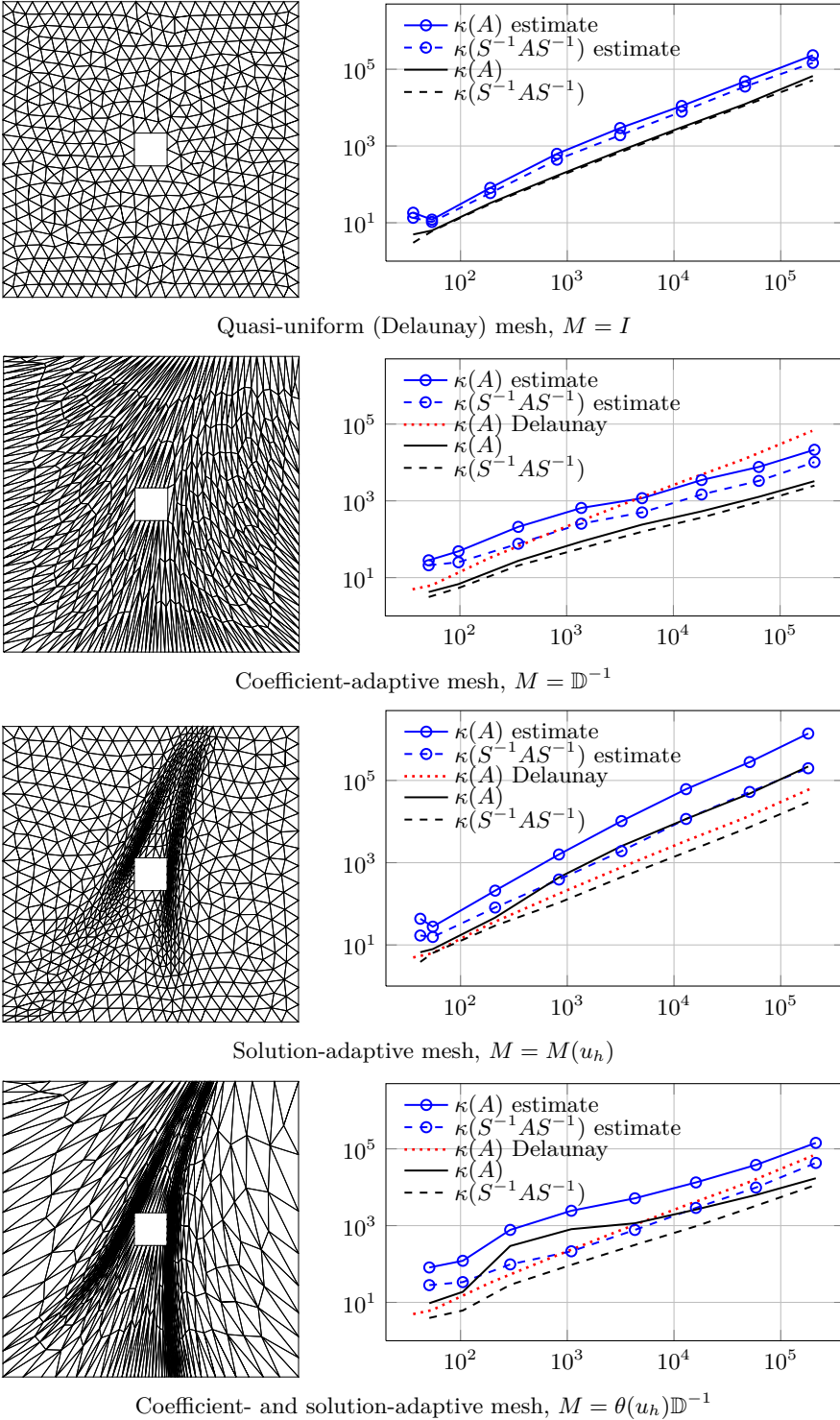


FIGURE 6. Example 6.4: condition number in dependence of  $N$

7. SUMMARY AND CONCLUSIONS

**Mass matrix.** Our new estimate (8) of the condition number of the Galerkin mass matrix is tight within a factor of  $(d + 2)$  from *both above and below* for *any mesh* with *no assumptions* on mesh regularity or topology. In this sense, it is optimal and truly anisotropic.

**Stiffness matrix.** Lemma 4.1 provides an estimate of the largest eigenvalue of the stiffness matrix which is simple to compute and is tight within a factor of  $(d + 1)$  from *both above and below* for *any mesh*. This is in contrast to many existing estimates which are proportional to the maximal number of elements meeting at a mesh point.

New bounds (21)–(23) on the smallest eigenvalue and (26)–(29) on the condition number of the stiffness matrix are a significant improvement in comparison to the previously available estimates.

First, the new bounds show that the conditioning of the stiffness matrix with an arbitrary (anisotropic) mesh is much better than generally assumed, especially for  $d = 1$  and  $d = 2$ .

Second, the new bounds are truly anisotropic and valid for any mesh since no assumptions on the mesh regularity were made.

Third, bounds (26) and (27) reveal what affects the conditioning. There are three factors. The first (base) factor  $CN^{\frac{2}{d}}$  describes the direct dependence of the condition number on the *number of mesh elements* and corresponds to the condition number for the Laplace operator on a uniform mesh. The second factor describes the effects of the *mesh  $\mathbb{D}^{-1}$ -nonuniformity*, i.e., the interplay between the shape and size of mesh elements and the coefficients of the BVP. It is  $\mathcal{O}(1)$  for a coefficient-adaptive mesh, i.e., a mesh satisfying (19). The third factor measures how the *mesh volume-nonuniformity* further affects the condition number. It has no effect in 1D, a minimal one in 2D, and a substantial effect in 3D and higher dimensions. This means that even if the mesh is coefficient-adaptive and the second factor is  $\mathcal{O}(1)$ , the mesh volume-nonuniformity can still have a significant impact on the condition number for  $d \geq 3$ .

Fourth, a simple diagonal scaling, such as the Jacobi preconditioning, can significantly improve the conditioning. Bound (29) for the condition number after scaling does not contain the factor for the mesh volume-nonuniformity. As a consequence, for a coefficient-adaptive mesh, this bound reduces to the base factor  $CN^{\frac{2}{d}}$ . In this sense, diagonal scaling eliminates the effects of the mesh volume-nonuniformity. It can also significantly reduce the effects of the mesh nonuniformity with respect to  $\mathbb{D}^{-1}$ : the influence reduces essentially from the maximum norm to the  $L^{\max\{1, \frac{d}{2}\}}$  norm of  $\|(F'_K)^{-1}\mathbb{D}_K(F'_K)^{-T}\|_2$ .

Moreover, for a preconditioner that is invariant to diagonal scaling it follows that the condition number of the preconditioned stiffness matrix is typically smaller than  $\kappa(S^{-1}AS^{-1})$  which in turn has a much lower bound than  $\kappa(A)$  (cf. (27) and (29)). For example, consider an incomplete Cholesky decomposition of  $A$ ,

$$A = LL^T + E.$$

It follows that

$$S^{-1}AS^{-1} = (S^{-1}L)(S^{-1}L)^T + S^{-1}ES^{-1}$$

is actually an incomplete Cholesky decomposition of  $S^{-1}AS^{-1}$  since  $S^{-1}L$  has the same sparsity pattern as  $L$ . Then from the identity

$$L^{-1}AL^{-T} = (S^{-1}L)^{-1} (S^{-1}AS^{-1}) (S^{-1}L)^{-T}$$

we see that the preconditioned matrix of  $A$  with preconditioner  $L$  is equivalent to the preconditioned matrix of  $S^{-1}AS^{-1}$  with preconditioner  $S^{-1}L$ . As a result, the performance of the preconditioning technique on  $A$  is the same as that on  $S^{-1}AS^{-1}$ , which has a much smaller condition number than  $A$ . Although there is no estimate yet on the condition number of the preconditioned system, the above observation may provide a partial explanation for the good performance of ILU preconditioners with anisotropic meshes observed in [12].

Numerical experiments (Figures 3a and 4a) indicate that although the new bounds have the same order as the exact value as the number of elements increases, they may have higher asymptotic orders than the exact value as the element aspect ratio increases. These may deserve further investigation.

Finally, we would like to point out that although the study in this paper has been done specifically for the linear finite element discretization, the approach can be generalized for higher order finite elements without major modifications.

#### ACKNOWLEDGEMENT

The first author is very thankful to Jonathan R. Shewchuk for a fruitful discussion at the ICIAM 2011 and for pointing out valuable references. The authors are grateful to the anonymous referees for their comments and suggestions which helped significantly to improve the quality of this paper.

#### REFERENCES

- [1] Thomas Apel, *Anisotropic finite elements: local estimates and applications*, Advances in Numerical Mathematics, B. G. Teubner, Stuttgart, 1999. MR1716824 (2000k:65002)
- [2] Randolph E. Bank and L. Ridgway Scott, *On the conditioning of finite element equations with highly refined meshes*, SIAM J. Numer. Anal. **26** (1989), no. 6, 1383–1394, DOI 10.1137/0726080. MR1025094 (90m:65192)
- [3] Susanne C. Brenner and L. Ridgway Scott, *The mathematical theory of finite element methods*, 3rd ed., Texts in Applied Mathematics, vol. 15, Springer, New York, 2008. MR2373954 (2008m:65001)
- [4] E. F. D’Azevedo, C. H. Romine, and J. M. Donato, *Coefficient adaptive triangulation for strongly anisotropic problems*, Tech. Report ORNL/TM-13086, Oak Ridge National Laboratory, 1997.
- [5] Qiang Du, Desheng Wang, and Liyong Zhu, *On mesh geometry and stiffness matrix conditioning for general finite element spaces*, SIAM J. Numer. Anal. **47** (2009), no. 2, 1421–1444, DOI 10.1137/080718486. MR2497335 (2010b:65252)
- [6] Alexandre Ern and Jean-Luc Guermond, *Theory and practice of finite elements*, Applied Mathematical Sciences, vol. 159, Springer-Verlag, New York, 2004. MR2050138 (2005d:65002)
- [7] Isaac Fried, *Bounds on the spectral and maximum norms of the finite element stiffness, flexibility and mass matrices* (English, with Russian summary), Internat. J. Solids and Structures **9** (1973), 1013–1034. MR0345400 (49 #10136)
- [8] David Gilbarg and Neil S. Trudinger, *Elliptic partial differential equations of second order*, Classics in Mathematics, Springer-Verlag, Berlin, 2001. Reprint of the 1998 edition. MR1814364 (2001k:35004)
- [9] F. Hecht, *Bamg: Bidimensional anisotropic mesh generator*, <http://www.ann.jussieu.fr/hecht/ftp/bamg/>.
- [10] Nicholas J. Higham, *Accuracy and stability of numerical algorithms*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996. MR1368629 (97a:65047)



- [11] Weizhang Huang, *Measuring mesh qualities and application to variational mesh adaptation*, SIAM J. Sci. Comput. **26** (2005), no. 5, 1643–1666 (electronic), DOI 10.1137/S1064827503429405. MR2142589 (2006b:65187)
- [12] Weizhang Huang, *Metric tensors for anisotropic mesh generation*, J. Comput. Phys. **204** (2005), no. 2, 633–665, DOI 10.1016/j.jcp.2004.10.024. MR2131856 (2005k:65274)
- [13] Weizhang Huang, *Mathematical principles of anisotropic mesh adaptation*, Commun. Comput. Phys. **1** (2006), no. 2, 276–310.
- [14] Weizhang Huang, Lennard Kamenski, and Jens Lang, *A new anisotropic mesh adaptation method based upon hierarchical a posteriori error estimates*, J. Comput. Phys. **229** (2010), no. 6, 2179–2198, DOI 10.1016/j.jcp.2009.11.029. MR2586243 (2011a:65406)
- [15] Weizhang Huang and Robert D. Russell, *Adaptive moving mesh methods*, Applied Mathematical Sciences, vol. 174, Springer, New York, 2011. MR2722625 (2012a:65243)
- [16] L. Kamenski, *A study on using hierarchical basis error estimates in anisotropic mesh adaptation for the finite element method*, Eng. Comput. **28** (2012), no. 4, 451–460.
- [17] Xianping Li and Weizhang Huang, *An anisotropic mesh adaptation method for the finite element solution of heterogeneous anisotropic diffusion problems*, J. Comput. Phys. **229** (2010), no. 21, 8072–8094, DOI 10.1016/j.jcp.2010.07.009. MR2719161 (2011i:65218)
- [18] J. R. Shewchuk, *What is a good linear finite element? Interpolation, conditioning, anisotropy, and quality measures*, <http://www.cs.cmu.edu/~jrs/jrspapers.html#quality>, 2002.
- [19] A. J. Wathen, *Realistic eigenvalue bounds for the Galerkin mass matrix*, IMA J. Numer. Anal. **7** (1987), no. 4, 449–457, DOI 10.1093/imanum/7.4.449. MR968517 (90a:65246)

DEPARTMENT OF MATHEMATICS, THE UNIVERSITY OF KANSAS, LAWRENCE, KANSAS 66045  
*Current address:* Weierstrass Institute, Mohrenstr. 39, 10117 Berlin, Germany  
*E-mail address:* kamenski@wias-berlin.de

DEPARTMENT OF MATHEMATICS, THE UNIVERSITY OF KANSAS, LAWRENCE, KANSAS 66045  
*E-mail address:* whuang@ku.edu

DEPARTMENT OF MATHEMATICS, THE UNIVERSITY OF KANSAS, LAWRENCE, KANSAS 66045  
*E-mail address:* xu@math.ku.edu