

An *hp* Finite Element Method for convection-diffusion problems

J.M. Melenk and C. Schwab

Research Report No. 97-05
February 1997

Seminar für Angewandte Mathematik
Eidgenössische Technische Hochschule
CH-8092 Zürich
Switzerland

An hp Finite Element Method for convection-diffusion problems

J.M. Melenk and C. Schwab

Seminar für Angewandte Mathematik
Eidgenössische Technische Hochschule
CH-8092 Zürich
Switzerland

Research Report No. 97-05

February 1997

Abstract

We analyze an hp FEM for convection-diffusion problems. Stability is achieved by suitably upwinded test functions, generalizing the classical α -quadratically upwinded and the Hemker test-functions for piecewise linear trial spaces (see, e.g., [12] and the references there). The method is proved to be stable independently of the viscosity. Further, the stability is shown to depend only weakly on the spectral order. We show how sufficiently accurate, approximate upwinded test functions can be computed on each element by a local least squares FEM. Under the assumption of analyticity of the input data, we prove robust exponential convergence of the method. Numerical experiments confirm our convergence estimates and show robust exponential convergence of the hp -FEM even for viscosities of the order of machine precision, i.e., for the limiting transport problem.

Subject Classification: Primary: 65N30, 65N12; Secondary: 35B25, 35C20

1 Convection-Diffusion Problem

1.1 Problem Formulation

In $\Omega = (-1, 1)$ we consider the model convection diffusion problem

$$L_\varepsilon u_\varepsilon := -\varepsilon u_\varepsilon'' + a(x)u_\varepsilon' + b(x)u_\varepsilon = f(x) \quad (1.1)$$

with the boundary conditions

$$u(\pm 1) = \alpha^\pm \in \mathbb{R}. \quad (1.2)$$

Here $\varepsilon > 0$ is the diffusivity, $u(x)$ is, for example, concentration of a transported substance, $a(x)$ is the velocity of the transporting medium, $b(x)$ specifies losses/ sources of the substance and $f(x)$ is an external source term. Throughout this work, we make the following assumptions on the coefficients $a \in C^1[-1, 1]$, $b \in C^0[-1, 1]$: There are constants $\underline{b} \in \mathbb{R}$, \underline{a} , $\gamma_1, \gamma_2 > 0$ such that for all $\varepsilon \in (0, 1]$

$$\begin{aligned} a(x) \geq \underline{a}, \quad b(x) \geq \underline{b}, \quad \min \{\underline{a}^2, \underline{a}^2 + 4\underline{b}\varepsilon\} \geq \gamma_1^2, \\ \max \left\{ \frac{\underline{a} - \sqrt{\underline{a}^2 + 4\underline{b}\varepsilon}}{2\varepsilon}, 0 \right\} \leq \gamma_2. \end{aligned} \quad (1.3)$$

As we will also consider the adjoint problem of (1.1) (cf. Section 2 ahead) we stipulate the existence of $\underline{b}^* \in \mathbb{R}$ and $\gamma_1^*, \gamma_2^* > 0$ such that

$$\begin{aligned} a(x) \geq \underline{a} > 0, \quad b(x) - a'(x) \geq \underline{b}^*, \quad \min \{\underline{a}^2, \underline{a}^2 + 4\underline{b}^*\varepsilon\} \geq (\gamma_1^*)^2, \\ \max \left\{ \frac{\underline{a} - \sqrt{\underline{a}^2 + 4\underline{b}^*\varepsilon}}{2\varepsilon}, 0 \right\} \leq \gamma_2^*. \end{aligned} \quad (1.4)$$

(1.3) ensures the unique solvability of (1.1), (1.2) while (1.4) guarantees the unique solvability of the adjoint problem. Note that for given a, b , the constants $\gamma_1, \gamma_2, \gamma_1^*, \gamma_2^*$ exist under the assumption that the diffusivity ε is sufficiently small.

The finite element approximation of (1.1), (1.2) for small ε is nontrivial due to the singular perturbation character of the problem which manifests itself in two distinct phenomena: First, the solution u_ε exhibits a boundary layer near the outflow boundary $x = 1$; we will characterize the boundary layer behavior of the solution u_ε more precisely in Section 1.2 ahead. The second difficulty arises from the well-known fact that symmetric variational formulations of (1.1), (1.2) based on H_0^1 as trial and test space are not uniformly stable with respect to the parameter ε . One possible remedy is the use of streamline-diffusion techniques which amount in effect to a nonconforming method (see, e.g., [14] and the references there). Crucial to the convergence analysis of streamline diffusion FEM are H^2 regularity of the solution and certain elementwise inverse inequalities which allow to control the higher order derivatives introduced into the variational formulation through the streamline-diffusion term. While this idea is, in principle, also feasible for p - or spectral element methods, the convergence rates obtained that way will be suboptimal due to

the higher loss of derivatives in inverse inequalities for polynomials. This is even more pronounced for the combined hp -type FEM, in particular in two and three dimensions, where the optimal approximation of boundary layers mandates elements of arbitrary high aspect ratio (see [16]) for which suitable inverse inequalities do not seem to be available. The hp -FEM is nevertheless very attractive for such problems, since hp -trial spaces have approximation properties superior to both, h -version FEM and spectral methods. For example, hp -FEM can be shown to approximate boundary layers and corner singularities at a *robust exponential convergence rate* (see [15, 16]). This will also be true for any stable projection method based on these trial spaces. To achieve stability, we use Petrov-Galerkin methods where test- and trial-spaces are distinct, an approach that has been followed by numerous authors in the finite difference and finite volume setting (see [12, 14, 11] for an account of this work and many references). We base the hp -FEM on the variational framework from [18, 19] where the h -version FEM was analyzed and optimal convergence rates, uniform in ε , were shown. Asymptotically exact h -version a-posteriori error estimators for this variational formulation have also been developed in [19] and it was shown that the numerical solutions exhibit few spurious oscillations and good pointwise convergence. The crucial ingredient in [18, 19] was the construction of suitable, upwinded test functions by asymptotic analysis of the elemental adjoint problem. The generalization of this asymptotic analysis to high order elements and higher dimensions is not straightforward.

Here we propose therefore a fully numerical method. More precisely, we show how for hp -trial spaces with any mesh-degree combination sufficiently accurate *approximate upwinded test functions* can be stably computed. The calculation of the test functions is completely localized to either a single element or a patch of elements and done by a least-squares like method (which is uniformly stable in ε). This can be simply performed as part of the usual element stiffness matrix generation in the hp -FEM. Our analysis shows that a) the approximate test functions thus obtained do ensure stability and that b) already fairly crude approximations of the test functions suffice, so that the work spent in computing these test functions can be expected to be moderate. Most importantly, no analytical input in the form of asymptotic expansions or boundary layer functions is necessary – the method is fully computational and conceptually generalizes readily to two- and three-dimensional problems. Here we analyze the method in detail for the one-dimensional model problem (1.1), (1.2), where new regularity results for the solution allow us to prove robust exponential convergence. The analysis for two and three dimensional problems will be given elsewhere [10].

1.2 Regularity

Let us consider (1.1) on $\Omega = (-1, 1)$ with *analytic* input data $a(x)$, $b(x)$, $f(x)$ satisfying

$$\|a^{(n)}\|_{L^\infty(\Omega)} \leq C_a \gamma_a^n \quad \forall n \in \mathbb{N}_0 \quad (1.5)$$

$$\|b^{(n)}\|_{L^\infty(\Omega)} \leq C_b \gamma_b^n \quad \forall n \in \mathbb{N}_0 \quad (1.6)$$

$$\|f^{(n)}\|_{L^\infty(\Omega)} \leq C_f \gamma_f^n \quad \forall n \in \mathbb{N}_0 \quad (1.7)$$

for some constants $C_a, C_b, C_f, \gamma_a, \gamma_b, \gamma_f > 0$. Assumptions (1.3) and (1.5)–(1.7) ensure the existence of a unique, analytic solution u_ε of (1.1), (1.2). The purpose of this subsection

is to illuminate the regularity properties of u_ε in dependence on the parameter ε and the constants of (1.3), (1.5)–(1.7). These regularity results are necessary for the proof of *robust exponential convergence* of the *hp*-FEM obtained in the present paper. Although regularity results related to the ones presented here are in the literature ([14], [12]), the specific derivative bounds seem to be new (see also [9] for the related case of a reaction diffusion equation).

The solution u_ε of (1.1), (1.2) is analytic on Ω ; however, for small values of ε , it exhibits a boundary layer at the outflow boundary. This boundary behavior can be characterized with the aid of asymptotic expansions: For any expansion order $M \in \mathbb{N}_0$, we have the standard decomposition (see, e.g., [4])

$$u_\varepsilon = w_M + C_M u_\varepsilon^+ + r_M. \quad (1.8)$$

Here, u_M is the *asymptotic part* given by

$$\begin{aligned} w_M &:= \sum_{j=0}^M \varepsilon^j u_j + \alpha^- e^{-\Lambda(x)} \\ u_{j+1}(x) &:= e^{-\Lambda(x)} \int_{-1}^x \frac{e^{\Lambda(t)}}{a(t)} u_j''(t) dt \quad j = 0, \dots, M-1 \\ u_0(x) &:= e^{-\Lambda(x)} \int_{-1}^x \frac{e^{\Lambda(t)}}{a(t)} f(t) dt \\ \Lambda(x) &:= \int_{-1}^x \lambda(t) dt \\ \lambda(x) &:= \frac{b(x)}{a(x)} \end{aligned}$$

The *outflow boundary layer* u_ε^+ solves the problem

$$L_\varepsilon u_\varepsilon^+ = 0 \quad \text{on } \Omega, \quad u_\varepsilon^+(-1) = 0, \quad u_\varepsilon^+(1) = 1, \quad (1.9)$$

and C_M is given by

$$C_M := \alpha^+ - w_M(1). \quad (1.10)$$

Finally, the remainder r_M is given as the solution of

$$L_\varepsilon r_M = \varepsilon^{M+1} u_M'' \quad \text{on } \Omega, \quad r_M(\pm 1) = 0 \quad (1.11)$$

Note that for $M = 0$ the function w_0 solves the *limiting transport problem* given by (1.1) with $\varepsilon = 0$ and the boundary condition $w_0(-1) = \alpha^-$.

Theorem 1.1 *Let u_ε be the solution of (1.1), (1.2). Then there are constants C, K depending only on the constants in (1.5)–(1.7) and on the constants $\underline{a}, \gamma_1, \gamma_2$ such that*

$$\|u_\varepsilon^{(n)}\|_{L^\infty(I)} \leq CK^n \max(n, \varepsilon^{-1})^n \quad \forall n \in \mathbb{N}_0 \quad (1.12)$$

$$|u_\varepsilon^{+(n)}(x)| \leq CK^n \max(n, \varepsilon^{-1})^n e^{-\underline{a}(1-x)/(2\varepsilon)} \quad \forall n \in \mathbb{N}_0, \quad x \in I. \quad (1.13)$$

Furthermore, under the assumption $0 < \varepsilon MK \leq 1$, the terms in the decomposition (1.8) satisfy

$$\|w_M^{(n)}\|_{L^\infty(I)} \leq CK^n n! \quad \forall n \in \mathbf{N}_0 \quad (1.14)$$

$$\|r_M^{(n)}\|_{L^\infty(I)} \leq C\varepsilon^{1-n}(\varepsilon MK)^M \quad n = 0, 1, 2 \quad (1.15)$$

$$|C_M| \leq C \quad (1.16)$$

The proof of Theorem 1.1 is given in Appendix B.

2 Variational Formulation

Without loss of generality, we may analyze (1.1) with homogeneous Dirichlet data

$$\alpha^\pm = 0 \quad (2.1)$$

by the standard argument of seeking u_ε in the form $u_\varepsilon = \tilde{u}_\varepsilon + u_0$ (where u_0 is linear and satisfies the boundary conditions (1.2)) and then noting that this leads to (1.1),(2.1) for \tilde{u}_ε with the same operator L_ε and suitably adjusted right hand side f which is analytic and independent of ε .

To motivate our variational formulation, we observe that multiplication of (1.1) by a test function v and twofold integration by parts gives a so-called *very weak variational formulation*: Find $u \in L^2(\Omega)$ such that

$$B(u, v) := \int_\Omega u L_\varepsilon^* v dx = \int_\Omega f v dx =: F(v) \quad \forall v \in H^2 \cap H_0^1(\Omega).$$

Here, L_ε^* denotes the adjoint of L_ε , i.e.

$$L_\varepsilon^* u = -\varepsilon u'' - a(x)u' + (b - a')(x)u \quad (2.2)$$

which is defined when $a \in C^1([-1, 1])$. There are several drawbacks with FEM based on very weak variational formulations: first, a' is in general not globally continuous, but only elementwise smooth (if it stems, for example, from linearization of the nonlinear problem around a FE-approximation of u), second, to obtain a good test-space for a given trial space of possibly discontinuous functions, a global adjoint problem must be solved for each basis function and third, the essential boundary conditions (1.2) are generally not satisfied by FE solutions. This leads us to a formulation which is situated “between” the weak one based on $H_0^1 \times H_0^1$ and the very weak one based on $L^2 \times H^2 \cap H_0^1$.

We present Sobolev spaces with mesh-dependent norms introduced in [19]. For a collection of nodes $\{-1 = x_0 < x_1 < \dots < x_N = 1\}$, we introduce the notation $I_j := (x_{j-1}, x_j)$, $h_j := |I_j| = x_j - x_{j-1}$, $m_j = (x_{j-1} + x_j)/2$ for $j = 1, \dots, N$. The elements I_j form a mesh $\mathcal{T} = \{I_j : j = 1, \dots, N\}$ on Ω . Let further $\{\rho_j\}_{j=1}^{N-1}$ be a sequence of positive numbers and set $\rho := \rho_1 + \rho_2 + \dots + \rho_{N-1}$, $h := \max\{h_j : j = 1, \dots, N\}$. Then we define the trial space $H_{\mathcal{T}}^0$ as completion of $H_0^1(\Omega)$ with respect to the mesh dependent norm

$$\|u\|_{H_{\mathcal{T}}^0} := \left(\int_{-1}^1 |u|^2 dx + \sum_{j=1}^{N-1} \rho_j |u(x_j)|^2 \right)^{1/2}. \quad (2.3)$$

The space $H_{\mathcal{T}}^0$ thus obtained is a Hilbert space and is isomorphic to $L^2(\Omega) \oplus \mathbb{R}^{N-1}$ so that every $u \in H_{\mathcal{T}}^0$ is of the form $u = (\tilde{u}, d_1, d_2, \dots, d_{N-1})$ and

$$\|u\|_{H_{\mathcal{T}}^0} = \left(\|\tilde{u}\|_{L^2}^2 + \sum_{j=1}^{N-1} \rho_j |d_j|^2 \right)^{1/2}. \quad (2.4)$$

If $u \in H_{\mathcal{T}}^0 \cap H^1$ then $\tilde{u} \in H^1$ and $d_j = \tilde{u}(x_j)$. If $u \in L^2(\Omega)$ is discontinuous at the nodes x_j but piecewise smooth, then $\tilde{u} = u$ and $d_j = (u(x_j^+) + u(x_j^-))/2$ where $u(x_j^\pm)$ denotes the left and right limits of u at x_j .

Next, we introduce the test space

$$H_{\mathcal{T}}^2 := \left\{ v \in H_0^1(\Omega) : v|_{I_j} \in H^2(I_j), j = 1, \dots, N. \right\} \quad (2.5)$$

On the pair $H_{\mathcal{T}}^0 \times H_{\mathcal{T}}^2$ we define the bilinear form $B_{\mathcal{T}}(\cdot, \cdot)$ by

$$B_{\mathcal{T}}(u, v) := \sum_{j=1}^N \int_{I_j} \tilde{u} L_{\varepsilon}^* v dx - \sum_{j=1}^{N-1} d_j [\varepsilon v'(x_j)] \quad (2.6)$$

where $[v'(x_j)]$ denotes the jump of v' at $x_j \in \mathcal{T}$. We equip the space $H_{\mathcal{T}}^2$ with the norm

$$\|v\|_{H_{\mathcal{T}}^2} := \left(\sum_{j=1}^N \|L_{\varepsilon}^* v\|_{L^2(I_j)}^2 + \sum_{j=1}^{N-1} \frac{|[\varepsilon v'(x_j)]|^2}{\rho_j} \right)^{1/2}. \quad (2.7)$$

We remark in passing that so far we have used $a(x) \in C^0([-1, 1]) \cap C^1(\bar{I}_j)$, $j = 1, \dots, N$, rather than $a \in C^1[-1, 1]$. With these definitions we have

Proposition 2.1 *For any mesh \mathcal{T} and any positive sequence $\{\rho_j\}_{j=1}^{N-1}$, the bilinear form $B_{\mathcal{T}}(\cdot, \cdot)$ satisfies*

$$|B_{\mathcal{T}}(u, v)| \leq \|u\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2} \quad \forall u \in H_{\mathcal{T}}^0, v \in H_{\mathcal{T}}^2, \quad (2.8)$$

$$\inf_{0 \neq v \in H_{\mathcal{T}}^2} \sup_{0 \neq u \in H_{\mathcal{T}}^0} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2}} \geq 1, \quad (2.9)$$

and

$$\forall 0 \neq u \in H_{\mathcal{T}}^0 : \sup_{v \in H_{\mathcal{T}}^2} B_{\mathcal{T}}(u, v) > 0. \quad (2.10)$$

Proof: The bound (2.8) follows directly from the definition of the norms and the Schwarz inequality.

To show (2.9), for given $v \in H_{\mathcal{T}}^2$, we select $u_v = (\tilde{u}, d_1, \dots, d_{N-1}) \in H_{\mathcal{T}}^0$ as follows:

$$\begin{aligned} \tilde{u}|_{I_j} &= \operatorname{sgn}(L_{\varepsilon}^* v|_{I_j}) |L_{\varepsilon}^* v|_{I_j}|, \quad j = 1, \dots, N, \\ d_j &= -\rho_j^{-1} [v'(x_j)] \quad j = 1, \dots, N-1. \end{aligned}$$

Then $\|u_v\|_{H_{\mathcal{T}}^0} \leq \|v\|_{H_{\mathcal{T}}^2}$ and $B_{\mathcal{T}}(u_v, v) = \|v\|_{H_{\mathcal{T}}^2}^2$, whence for every $0 \neq v \in H_{\mathcal{T}}^2$

$$\sup_{0 \neq u \in H_{\mathcal{T}}^0} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2}} \geq \frac{B_{\mathcal{T}}(u_v, v)}{\|u_v\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2}} \geq \frac{B_{\mathcal{T}}(u_v, v)}{\|v\|_{H_{\mathcal{T}}^2}^2} = 1$$

which proves (2.9) and (2.10). □

Proposition 2.1 allows us to prove the existence of a solution $u \in H_{\mathcal{T}}^0$ of the problem:

$$u \in H_{\mathcal{T}}^0 \quad B_{\mathcal{T}}(u, v) = F(v) \quad \forall v \in H_{\mathcal{T}}^2. \quad (2.11)$$

Proposition 2.2 *Under the assumption (1.3), for every $f \in L^1(\Omega)$, every $0 < \varepsilon \leq 1$, and every mesh \mathcal{T} and every positive sequence $\{\rho_j\}_{j=1}^{N-1}$, the problem (2.11) admits a unique solution $u \in H_{\mathcal{T}}^0$.*

Proof: We proceed in two steps.

Step i) We claim that for any mesh \mathcal{T} , any positive sequence $\{\rho_j\}_{j=1}^{N-1}$ and any $\varepsilon \in (0, 1]$, assumptions (1.3) imply

$$\|v\|_{L^\infty(\Omega)} \leq C_1 \max\{1, \sqrt{\rho}\} \|v\|_{H_{\mathcal{T}}^2}, \quad (2.12)$$

where C_1 is independent of ε and \mathcal{T} . To prove it, we note that (1.3) implies the existence of a Green's function $G(x, y)$ for the problem (1.1), (1.2) which is bounded uniformly with respect to x, y, ε (see [18], Theorem 2.7), i.e.

$$\max_{(x,y) \in [-1,1]^2} |G(x, y)| \leq C_G.$$

For $v \in H_{\mathcal{T}}^2$, we can write

$$v(y) = \sum_{j=1}^N \int_{I_j} G(x, y) (L_\varepsilon^* v)(x) dx - \sum_{j=1}^{N-1} [\varepsilon v'(x_j)] G(x_j, y), \quad \forall y \in [-1, 1].$$

Using the boundedness of $G(x, y)$, we estimate then

$$\begin{aligned} |v(y)| &\leq C_G \left\{ \sum_{j=1}^N \int_{I_j} |L_\varepsilon^* v| dx + \sum_{j=1}^{N-1} \rho_j^{-1/2} |[\varepsilon v'(x_j)]| \rho_j^{1/2} \right\} \\ &\leq \sqrt{2} C_G \max\{\sqrt{2}, \sqrt{\rho}\} \|v\|_{H_{\mathcal{T}}^2} \end{aligned}$$

which proves (2.12).

Step ii) For $f \in L^1(\Omega)$ and $v \in H_{\mathcal{T}}^2$, we therefore have

$$|F(v)| \leq \|f\|_{L^1(\Omega)} \|v\|_{L^\infty(\Omega)} \leq C_1 \max\{1, \sqrt{\rho}\} \|f\|_{L^1(\Omega)} \|v\|_{H_{\mathcal{T}}^2}.$$

Hence, $F(\cdot)$ is a continuous, linear functional on $H_{\mathcal{T}}^2$ the norm of which is bounded uniformly with respect to ε and \mathcal{T} . By Propositions 2.1 and A.2, we have also

$$\begin{aligned} \inf_{0 \neq u \in H_{\mathcal{T}}^0} \sup_{0 \neq v \in H_{\mathcal{T}}^2} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2}} &\geq 1, \\ \forall 0 \neq v \in H_{\mathcal{T}}^2 : \sup_{u \in H_{\mathcal{T}}^0} B_{\mathcal{T}}(u, v) &> 0. \end{aligned}$$

This implies with $F \in (H_{\vec{\tau}}^2)'$ and Proposition A.1 that (2.11) admits a unique solution and that the a-priori estimate

$$\|u\|_{H_{\vec{\tau}}^0} \leq C_1 \max\{1, \sqrt{\rho}\} \|f\|_{L^1(\Omega)} \quad (2.13)$$

holds. □

Remark 2.3 In Step ii) of the proof of Proposition 2.2, we exploit the fact that the embedding $H_{\vec{\tau}}^2 \subset H_0^1(\Omega) \subset C(\overline{\Omega})$ is continuous. Hence, the right hand side functional F may actually be represented by an L^1 function f plus a (finite) number of Dirac distributions.

The variational formulation (2.11) is the basis of the FE-discretization.

3 hp Finite Element Discretization

3.1 The Finite Element spaces

We associate with each element I_j a polynomial degree $p_j \geq 1$ and combine the p_j in the degree-vector \vec{p} . We also set $p := \max\{p_j : 1 \leq j \leq N\}$. The trial spaces $S_0^{\vec{p}, \ell}(\mathcal{T})$ of our finite element method are the usual spaces of continuous, piecewise polynomials of degree p_j satisfying the homogeneous boundary conditions (2.1) at ± 1 :

$$\begin{aligned} S^{\vec{p}, \ell}(\mathcal{T}) &:= \left\{ u \in H^\ell(\Omega) : u|_{I_j} \in \Pi_{p_j}(I_j), j = 1, \dots, N \right\}, \quad \ell = 1, 2, \dots \\ S_0^{\vec{p}, \ell}(\mathcal{T}) &:= S^{\vec{p}, \ell}(\mathcal{T}) \cap H_0^1(\Omega) \end{aligned} \quad (3.1)$$

If $\ell = 1$, we simply write $S_0^{\vec{p}}(\mathcal{T})$.

As test space we choose, following [19], the space of L -splines of degree \vec{p} defined by

$$S_L^{\vec{p}}(\mathcal{T}) := \left\{ v \in H_0^1(\Omega) : (L_\varepsilon^* v)|_{I_j} = 0 \text{ if } p_j = 1, (L_\varepsilon^* v)|_{I_j} \in \Pi_{p_j-2}(I_j) \text{ if } p_j \geq 2 \right\}. \quad (3.2)$$

Note that (1.3), (1.4) imply that L_ε and L_ε^* are injective. Hence (3.2) makes sense and the test functions belong to $H^2(I_j)$, $j = 1, \dots, N$. We omit the argument \mathcal{T} when it is clear from the context which mesh is meant. Note that, due to (2.2), the space $S_L^{\vec{p}}$ is well-defined even if the coefficient $a(x)$ is only piecewise C^1 . We also observe that

$$M = \dim(S_0^{\vec{p}}) = -1 + \sum_{j=1}^N p_j = \dim(S_L^{\vec{p}}). \quad (3.3)$$

The finite element approximation u_M is then obtained in the usual way:

$$u_M \in S_0^{\vec{p}} \quad B_{\mathcal{T}}(u_M, v) = F(v) \quad \forall v \in S_L^{\vec{p}}. \quad (3.4)$$

Due to (3.3), problem (3.4) amounts to solving a linear system of M equations for the M unknown coefficients of u_M .

3.2 Stability

Our main result in this section is

Theorem 3.1 *Select*

$$\rho_j := (h_j + h_{j+1})/2, \quad j = 1, \dots, N-1. \quad (3.5)$$

Then for all $0 < \varepsilon \leq 1$, \mathcal{T} and \vec{p} there holds

$$\inf_{0 \neq v \in S_L^{\vec{p}}} \sup_{0 \neq u \in S_0^{\vec{p}}} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_{\mathcal{T}}^0} \|v\|_{H_{\mathcal{T}}^2}} \geq \frac{1}{\gamma_M} \quad (3.6)$$

with $\gamma_M = \max\{\sqrt{5}, \sqrt{p+3}\}$.

Proof: We show that for every $v \in S_L^{\vec{p}}$ there exists $u_v \in S_0^{\vec{p}}$ such that

$$B_{\mathcal{T}}(u_v, v) \geq \|v\|_{H_{\mathcal{T}}^2}^2, \quad \|u_v\|_{H_{\mathcal{T}}^0} \leq \gamma_M \|v\|_{H_{\mathcal{T}}^2}.$$

To this end, we write

$$u_v|_{I_j} = \sum_{i=0}^{p_j} a_{ij} L_i \left(2 \frac{x - m_j}{h_j} \right), \quad (3.7)$$

where L_i denotes the i th Legendre polynomial on $(-1, 1)$ normalized such that $L_i(1) = 1$. A basis for $S_L^{\vec{p}}$ can be obtained as follows: First, we define *external, nodal upwinded shape functions* $\psi_{-1,j} \in H_0^1(\bar{I}_{j-1} \cup \bar{I}_j)$ by

$$\begin{aligned} L_{\varepsilon}^* \psi_{-1,j} &= 0 \text{ in } I_{j-1} \cup I_j, j = 2, \dots, N, \\ \psi_{-1,j}(x_k) &= \delta_{j,k+1}, k = 1, \dots, N, \\ \psi_{-1,j} &= 0 \text{ elsewhere.} \end{aligned} \quad (3.8)$$

Note that $\psi_{-1,j} \in H^2(I_j)$, $j = 1, \dots, N$. The nodal shape functions $\psi_{-1,j}$ are augmented for $p_j \geq 2$ by *internal, upwinded shape functions* $\psi_{i,j} \in (H^2 \cap H_0^1)(I_j)$. They are defined by

$$\begin{aligned} L_{\varepsilon}^* \psi_{i,j} &= L_i \left(2 \frac{x - m_j}{h_j} \right) \text{ in } I_j, i = 0, \dots, p_j - 2, j = 1, \dots, N \\ \psi_{i,j} &= 0 \text{ elsewhere.} \end{aligned} \quad (3.9)$$

Any $v \in S_L^{\vec{p}}$ can be written as

$$v(x) = \sum_{j=2}^N v(x_{j-1}) \psi_{-1,j}(x) + \sum_{j=1}^N \sum_{i=0}^{p_j-2} b_{ij} \psi_{i,j}(x) \quad (3.10)$$

where b_{ij} are the Legendre coefficients of $L_{\varepsilon}^* v|_{I_j}$. Further, from the definition (3.8) of the $\psi_{-1,j}$ we have

$$L_{\varepsilon}^* v|_{I_j} = \sum_{i=0}^{p_j-2} b_{ij} L_i \left(2 \frac{x - m_j}{h_j} \right), \quad j = 1, \dots, N$$

which yields with the orthogonality properties of the Legendre polynomials and a scaling argument

$$\sum_{j=1}^N \|L_\varepsilon^* v\|_{L^2(I_j)}^2 = \sum_{j=1}^N h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} =: \sum_{j=1}^N h_j S_j. \quad (3.11)$$

Combining (3.11) with (2.7), we obtain for $v \in S_L^{\vec{p}}$ an expression for $\|v\|_{H_\tau^2}$ in terms of the b_{ij} :

$$\|v\|_{H_\tau^2} = \left(\sum_{j=1}^N h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} + \sum_{j=1}^{N-1} \frac{[\varepsilon v'(x_j)]^2}{\rho_j} \right)^{1/2} \quad \forall v \in S_L^{\vec{p}}. \quad (3.12)$$

Writing $v(x)$ in the form (3.10) and $u_v(x)$ as in (3.7) and inserting into (2.6), we find in the same way

$$B_{\mathcal{T}}(u_v, v) = \sum_{j=1}^{N-1} \left\{ h_j \sum_{i=0}^{p_j-2} \frac{a_{ij} b_{ij}}{2i+1} - \sum_{j=1}^{N-1} u_v(x_j) [\varepsilon v'(x_j)] \right\}. \quad (3.13)$$

For given $v \in S_L^{\vec{p}}$, i.e. for given b_{ij} , we choose now a_{ij} as follows: first, we select

$$a_{ij} = b_{ij} \quad i = 0, \dots, p_j - 2 \quad (3.14)$$

which leaves $a_{p_j-1,j}, a_{p_j,j}$ to be determined, for each I_j . Since u_v must be continuous, two conditions per interval must be enforced. We prescribe u_v at each endpoint of I_j as follows (by $u_v(x^\pm)$ we denote the right/left limit of u_v at x):

$$u_v(x_{j-1}^+) = a_j^- := \begin{cases} -[\varepsilon v'(x_{j-1})] / \rho_{j-1} & \text{if } j > 1, \\ 0 & \text{if } j = 1 \end{cases} \quad (3.15)$$

and

$$u_v(x_j^-) = a_j^+ := \begin{cases} -[\varepsilon v'(x_j)] / \rho_j & \text{if } j < N, \\ 0 & \text{if } j = N. \end{cases} \quad (3.16)$$

Conditions (3.15), (3.16) ensure continuity of u_v . Since $L_i(\pm 1) = (\pm 1)^i$ implies

$$u_v(x_{j-1}^+) = \sum_{i=0}^{p_j} (-1)^i a_{ij}, \quad u_v(x_j^-) = \sum_{i=0}^{p_j} (-1)^i a_{ij}$$

we get with (3.14) the linear system

$$\begin{bmatrix} (-1)^{p_j-1} & (-1)^{p_j} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} a_{p_j-1,j} \\ a_{p_j,j} \end{bmatrix} = \begin{bmatrix} a_j^- - \sum_{i=0}^{p_j-2} (-1)^i b_{ij} \\ a_j^+ - \sum_{i=0}^{p_j-2} b_{ij} \end{bmatrix}. \quad (3.17)$$

Its determinant is nonzero for any p_j , therefore u_v is uniquely determined by (3.14) and (3.17).

From (3.13), (3.14) and (3.15), (3.16) we get

$$B_{\mathcal{T}}(u_v, v) = \|v\|_{H_\tau^2}^2.$$

It remains therefore to show

$$\|u_v\|_{H_T^0} \leq \gamma_M \|v\|_{H_T^2} \quad (3.18)$$

with γ_M as in (3.6).

Since u_v is continuous, we have

$$\begin{aligned} \|u_v\|_{H_T^0}^2 &= \sum_{j=1}^N \|u_v\|_{L^2(I_j)}^2 + \sum_{j=1}^{N-1} \rho_j |u_v(x_j)|^2 \\ &= \sum_{j=1}^N h_j \sum_{i=0}^{p_j} \frac{|a_{ij}|^2}{2i+1} + \sum_{j=1}^{N-1} \rho_j |a_j^+|^2 \\ &= \sum_{j=1}^N \left\{ h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} + h_j \sum_{i=p_j-1}^{p_j} \frac{|a_{ij}|^2}{2i+1} \right\} + \sum_{j=1}^{N-1} \rho_j^{-1} |[\varepsilon v'(x_j)]|^2 \\ &= \|v\|_{H_T^2}^2 + \sum_{j=1}^N h_j \sum_{i=p_j-1}^{p_j} \frac{|a_{ij}|^2}{2i+1}. \end{aligned} \quad (3.19)$$

We estimate $|a_{ij}|^2$ for $i = p_j - 1, p_j$. From (3.17), we get

$$\begin{aligned} \begin{bmatrix} a_{p_j-1,j} \\ a_{p_j,j} \end{bmatrix} &= \frac{1}{2(-1)^{p_j}} \begin{pmatrix} 1 & (-1)^{p_j-1} \\ -1 & (-1)^{p_j-1} \end{pmatrix} \begin{pmatrix} a_j^- - \sum_{i=0}^{p_j-2} (-1)^i b_{ij} \\ a_j^+ - \sum_{i=0}^{p_j-2} b_{ij} \end{pmatrix} \\ &= \frac{1}{2(-1)^{p_j}} \begin{pmatrix} a_j^- + (-1)^{p_j-1} a_j^+ + \sum_{i=0}^{p_j-2} b_{ij} ((-1)^{p_j} - (-1)^i) \\ -a_j^- + (-1)^{p_j-1} a_j^+ + \sum_{i=0}^{p_j-2} b_{ij} ((-1)^i + (-1)^{p_j}) \end{pmatrix}. \end{aligned}$$

We estimate

$$\max\{|a_{ij}| : i = p_j - 1, p_j\} \leq \frac{1}{2} (|a_j^-| + |a_j^+|) + \sum_{i=0}^{p_j-2} |b_{ij}|$$

and get with (3.15), (3.16) that

$$\max\{|a_{ij}|^2 : i = p_j - 1, p_j\} \leq \varepsilon^2 \left(\frac{|[v'(x_{j-1})]|^2}{\rho_{j-1}^2} + \frac{|[v'(x_j)]|^2}{\rho_j^2} \right) + 2 \left(\sum_{i=0}^{p_j-2} |b_{ij}| \right)^2.$$

With the understanding that $[v'(x_0)] = [v'(x_N)] = 0$ and $\rho_0 = \rho_N = \infty$ we estimate further

$$\begin{aligned} &\sum_{j=1}^N h_j \sum_{i=p_j-1}^{p_j} \frac{|a_{ij}|^2}{2i+1} \\ &\leq \sum_{j=1}^N \frac{h_j}{2p_j-1} \left\{ \frac{|[\varepsilon v'(x_{j-1})]|^2}{\rho_{j-1}^2} + \frac{|[\varepsilon v'(x_j)]|^2}{\rho_j^2} + 2 \left(\sum_{i=0}^{p_j-2} |b_{ij}| \right)^2 \right\} \\ &\leq \sum_{j=1}^N \frac{h_j}{2p_j-1} \left\{ \frac{|[\varepsilon v'(x_{j-1})]|^2}{\rho_{j-1}^2} + \frac{|[\varepsilon v'(x_j)]|^2}{\rho_j^2} + 2 \sum_{i=0}^{p_j-2} (2i+1) \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} \right\} \\ &= \sum_{j=1}^N \frac{h_j}{2p_j-1} \left\{ \frac{|[\varepsilon v'(x_{j-1})]|^2}{\rho_{j-1}^2} + \frac{|[\varepsilon v'(x_j)]|^2}{\rho_j^2} + 2(p_j-1)^2 \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} \right\}. \end{aligned}$$

Now using $h_j/\rho_j \leq 2$, $h_j/\rho_{j-1} \leq 2$ and

$$\max\left\{\frac{2(p_j - 1)^2}{2p_j - 1} : j = 1, \dots, N\right\} \leq p + 2$$

we arrive at

$$\begin{aligned} \sum_{j=1}^N h_j \sum_{i=p_j-1}^{p_j} \frac{|a_{ij}|^2}{2i+1} &\leq 4 \sum_{j=1}^{N-1} \frac{|\varepsilon v'(x_j)|^2}{\rho_j} + (p+2) \sum_{j=1}^N h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} \\ &\leq \max\{4, p+2\} \|v\|_{H_{\mathcal{T}}}^2 \end{aligned}$$

where we used (2.7) and (3.11). Referring to (3.19) completes the proof. \square

Remark 3.2 In Theorem 3.1, we selected a specific sequence $\{\rho_j\}$. Inspection of the proof shows that any positive sequence is admissible. Then, however,

$$\gamma_M \geq C\sqrt{p} \max_{1 \leq j \leq N} \{h_j/\rho_j, h_j/\rho_{j-1}\}. \quad (3.20)$$

This shows that in order for γ_M to be independent of \mathcal{T} , the weights ρ_j must essentially be of the order of the local meshwidth.

3.3 Consistency and Convergence

Theorem 3.1 implies with Proposition A.4 and (2.8), (2.9) that

$$\inf_{0 \neq u \in S_0^{\vec{p}}} \sup_{0 \neq v \in S_L^{\vec{p}}} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_{\mathcal{T}}} \|\|v\|_{H_{\mathcal{T}}}^2} \geq \frac{1}{\gamma_M}, \quad (3.21)$$

$$\forall 0 \neq v \in S_L^{\vec{p}} : \sup_{u \in S_0^{\vec{p}}} B(u, v) > 0. \quad (3.22)$$

Referring to (2.10), we deduce from (3.21) and from Proposition A.3 that for every mesh-degree combination (\vec{p}, \mathcal{T}) there exists a unique FE-solution u_M of (3.4). In particular, the $M \times M$ (generally nonsymmetric) stiffness matrix corresponding to (3.4) is nonsingular. Moreover, the FE-solution u_M is quasi-optimal, i.e.

$$\|u - u_M\|_{H_{\mathcal{T}}} \leq (1 + \gamma_M) \|u - w\|_{H_{\mathcal{T}}} \quad \forall w \in S_0^{\vec{p}}. \quad (3.23)$$

The rate of convergence of the FEM (3.4) is therefore determined by the approximability of the exact solution u from the trial space $S_0^{\vec{p}}$. We show that proper selection of the mesh \mathcal{T} and of the polynomial degree distribution \vec{p} yields an *exponential rate of convergence*, uniform in ε . We will consider the approximation of two types of solutions u_ε . In Section 3.3.1, we consider the case of analytic right hand sides f . In that case, the solution u_ε exhibits only a boundary layer at the outflow boundary and thus a “two-element” mesh with one small element in the outflow boundary layer. In Section 3.3.2, we analyze the approximation of solutions stemming from right hand sides that contain Dirac distributions (Note that such right hand sides are admissible by Remark 2.3). Such solutions exhibit internal layers and in the limit (as $\varepsilon \rightarrow 0$) such solutions correspond to shocks. We show in Section 3.3.2 that the addition of a small element in each resolves these smeared-out shocks robustly.

3.3.1 Approximation of boundary layers

Theorem 3.3 *Let u_ε be the solution of (1.1), (2.1), and assume that the coefficients a , b and the right hand side f satisfy (1.3), (1.5)–(1.7). For every ε , $\kappa > 0$ let the degree vector \vec{p} and the mesh $\mathcal{T} = \mathcal{T}_{\kappa,\varepsilon}$ be given by*

$$\begin{aligned} \vec{p} &= \{p, p\}, & \mathcal{T}_{\kappa,\varepsilon} &= \{I_1, I_2\} & \text{if } \kappa p \varepsilon < 1, \\ \vec{p} &= \{p\}, & \mathcal{T}_{\kappa,\varepsilon} &= \{\Omega\} & \text{if } \kappa p \varepsilon \geq 1. \end{aligned} \quad (3.24)$$

where

$$I_1 = (-1, 1 - \kappa p \varepsilon), \quad I_2 = (1 - \kappa p \varepsilon, 1).$$

Then there is a constant κ_0 depending only on the constants of (1.3), (1.5)–(1.7) such that for every $0 < \kappa < \kappa_0$ there are C , $\sigma > 0$ independent of p and ε such that

$$\inf_{v_p \in S_0^{\vec{p},1}(\mathcal{T}_{\kappa,\varepsilon})} \|u_\varepsilon - v_p\|_{L^\infty(\Omega)} \leq C e^{-\sigma p} \quad \forall p \in \mathbb{N}. \quad (3.25)$$

Let us comment on Theorem 3.3 before proving it. As $\|u_\varepsilon - v_p\|_{H_T^0}^2 \leq 3\|u_\varepsilon - v_p\|_{L^\infty(\Omega)}^2$, Theorem 3.3 shows with (3.23) that for analytic input data robust exponential convergence can be achieved by the FE scheme (3.4) provided the space $S_0^{\vec{p}}$ is designed properly (ie. with one element of size $O(p\varepsilon)$ in the outflow boundary layer) and provided that the corresponding stable test space $S_L^{\vec{p}}(\mathcal{T})$ is available. Results analogous to Theorem 3.3 hold also true when $f(x)$ is piecewise analytic on $[-1, 1]$; then, however, additional internal layers arise at points of nonanalyticity of f which must be accounted for by adding further $O(\varepsilon p)$ elements.

Remark 3.4 An estimate on the value of the constant κ_0 is in principle available from the proof of Theorem 3.3. For *constant* coefficients a , b , the value of κ_0 can be determined explicitly ([15]): $\kappa_0 = 4/(e\lambda)$ where $\lambda = (a + \sqrt{a^2 + 4b\varepsilon})/2 \geq a/2$ by assumption (1.3). The numerical experiments of [15] show moreover, that the approximation properties of piecewise polynomials on the meshes $\mathcal{T}_{\kappa,\varepsilon}$ are fairly insensitive to the precise choice of κ .

In order to prove Theorem 3.3, we need two lemmas on the approximation of analytic functions by their Gauss-Lobatto interpolants. Let $I = [-1, 1]$ and define on $C(I)$ interpolation operator i_p by interpolation in the $p + 1$ Gauss-Lobatto points. By [17] we have the following stability result

$$\|i_p u\|_{L^\infty(I)} \leq C_{GL}(1 + \ln p) \|u\|_{L^\infty(I)} \quad \forall u \in C(I), p \in \mathbb{N}. \quad (3.26)$$

A direct consequence of this stability estimate and Markov's inequality, $\|v_p'\|_{L^\infty(I)} \leq p^2 \|v_p\|_{L^\infty(I)}$, valid for all polynomials v_p of degree p , is the following

Lemma 3.5 *Let $u \in C^2(I)$. Then*

$$\|(u - i_p u)^{(l)}\|_{L^\infty(I)} \leq \|u^{(l)}\|_{L^\infty(I)} + C_{GL}(1 + \ln p) p^{2l} \|u\|_{L^\infty(I)}, \quad l = 0, 1, 2.$$

For analytic functions u , we have

Lemma 3.6 *Let $u \in C^\infty(I)$ satisfy*

$$\|u^{(n)}\|_{L^\infty(I)} \leq C_u \gamma^n n! \quad \forall n \in \mathbf{N}_0.$$

Then there are constants $C, \sigma > 0$ depending only on C_{GL} and γ such that

$$\|u - i_p u\|_{L^\infty(I)} + \|(u - i_p u)'\|_{L^\infty(I)} + \|(u - i_p u)''\|_{L^\infty(I)} \leq C C_u e^{-\sigma p} \quad \forall p \in \mathbf{N}.$$

Proof: The growth estimates on the derivatives of u imply that u is analytic on the closed set I . By standard theory, there are polynomials P_p of degree p (e.g., by interpolating u in the Tschebyscheff points; cf. [3], Chap. 4 for the details) such that

$$\|u - P_p\|_{L^\infty(I)} + \|(u - P_p)'\|_{L^\infty(I)} + \|(u - P_p)''\|_{L^\infty(I)} \leq C C_u e^{-\sigma p} \quad \forall p \in \mathbf{N}$$

for some $C, \sigma > 0$ depending only on γ . As the interpolation operator i_p reproduces polynomials, we have $u - i_p u = (u - P_p) - i_p(u - P_p)$ and the desired result follows from an application of Lemma 3.5 to $u - P_p$. □

Proof of Theorem 3.3: We will choose the approximant v_p as the (piecewise) Gauss-Lobatto interpolant of u_ε . Because the endpoints of the elements are sampling points of the Gauss-Lobatto interpolation operator and because $u_\varepsilon(\pm 1) = 0$, the piecewise Gauss-Lobatto interpolant is in $S_0^{\bar{p},1}(\mathcal{T}_{\kappa,\varepsilon})$, and we merely have to control the approximation error on the sub-intervals.

Let us first consider the asymptotic case, i.e., $\kappa p \varepsilon \geq 1$. By Theorem 1.1, we have

$$\|u_\varepsilon^{(n)}\|_{L^\infty(\Omega)} \leq C K^n \max(n, \varepsilon^{-1})^n \quad \forall n \in \mathbf{N}_0.$$

Furthermore, we have by Stirling's formula the existence of $C > 0$ such that

$$\max(n, \varepsilon^{-1})^n \leq \max(n^n, n! \varepsilon^{-n} / n!) \leq \max(n^n, n! e^{1/\varepsilon}) \leq C n! e^{1/\varepsilon}. \quad (3.27)$$

Lemma 3.6 allows us to conclude that there are $C, \sigma > 0$ independent of p, ε such that

$$\|u_\varepsilon - i_p u_\varepsilon\|_{L^\infty(\Omega)} \leq C e^{1/\varepsilon} e^{-\sigma p}.$$

The assumption $\kappa p \varepsilon \geq 1$ implies $e^{1/\varepsilon} \leq e^{\kappa p}$ and thus the claim of the theorem follows in the asymptotic regime provided that $0 < \kappa < \kappa_0 \leq \sigma$.

In the pre-asymptotic case $\kappa p \varepsilon < 1$, we use the decomposition (1.8) with expansion order M given by

$$M = \mu \kappa p \quad \text{with } \mu \text{ such that } \mu K =: \beta < 1 \quad (3.28)$$

where $K > 0$ is the constant of Theorem 1.1 (strictly speaking, we should take M as the integer part of $\mu \kappa p$ — for notational convenience, however, we will ignore this point henceforth). This choice of μ guarantees that the statements of Theorem 1.1 on the terms of the decomposition (1.8) hold true because $\kappa p \varepsilon \leq 1$. Denote by l_1 , and l_2 the two linear maps which map the reference interval I onto the physical elements I_1, I_2 , and define, for $u \in C[-1, 1]$, the piecewise Gauss-Lobatto interpolant $\pi_p(u) \in S^{\bar{p},1}(\mathcal{T}_{\kappa,\varepsilon})$ of u by

$$\pi_p(u)|_{I_i} = i_p(u \circ l_i) \circ l_i^{-1}. \quad (3.29)$$

Let us now consider the difference between the terms of the decomposition (1.8) and their Gauss-Lobatto interpolants.

First, let us analyze the term w_M . As the maps l_i are linear with $|l'_i| \leq 1$, $i = 1, 2$, Theorem 1.1 allows us to infer that the functions $w_M \circ l_i$ defined on the reference element I satisfy the derivative growth estimates

$$\| (w_M \circ l_i)^{(n)} \|_{L^\infty(I)} \leq CK^n n! \quad \forall n \in \mathbb{N}_0$$

with C, K given by Theorem 1.1. Thus, by Lemma 3.6 there are $C, \sigma > 0$ such that

$$\| w_M \circ l_i - i_p(w_M \circ l_i) \|_{L^\infty(I)} \leq Ce^{-\sigma p} \quad \forall p \in \mathbb{N}, \quad i = 1, 2$$

from whence we immediately get that

$$\| w_M - \pi_p(w_M) \|_{L^\infty(\Omega)} \leq Ce^{-\sigma p}.$$

Consider now the approximation of $C_M u_\varepsilon^+$ on the small boundary layer element I_2 . First, observe that $C_M \leq C$ independently of M, ε by our choice of μ . With the aid of Theorem 1.1, the fact that $l'_2 = \kappa p \varepsilon / 2$, and the assumption $\kappa p \varepsilon \leq 1$, we obtain

$$\begin{aligned} \| (u_\varepsilon^+ \circ l_2)^{(n)} \|_{L^\infty(I)} &\leq CK^n (\kappa p \varepsilon / 2)^n \max(n, \varepsilon^{-1})^n \leq C(K/2)^n \max(\kappa p \varepsilon n, \kappa p)^n \\ &\leq C(K/2)^n \max(n^n, n! (\kappa p)^n / n!) \leq C(K/2)^n n! e^n e^{\kappa p}. \end{aligned}$$

Hence, Lemma 3.6 allows us to conclude the existence of $C, \sigma > 0$ independent of ε, p such that

$$\| u_\varepsilon^+ \circ l_2 - i_p(u_\varepsilon^+ \circ l_2) \|_{L^\infty(I)} \leq Ce^{\kappa p} e^{-\sigma p}. \quad (3.30)$$

This term is exponentially small provided that $\kappa < \kappa_0 \leq \sigma$. Let us now turn our attention to the approximation of $C_M u_\varepsilon^+$ on I_1 . By Theorem 1.1, we have that

$$\| u_\varepsilon^+ \circ l_1 \|_{L^\infty(I)} \leq Ce^{-\underline{a}\kappa p / 2}.$$

Thus, by Lemma 3.5

$$\| u_\varepsilon^+ \circ l_1 - i_p(u_\varepsilon^+ \circ l_1) \|_{L^\infty(I)} \leq Ce^{-\sigma p} \quad (3.31)$$

for some properly chosen $\sigma > 0$. Combining (3.30), (3.31) allows us to conclude that

$$\| C_M u_\varepsilon^+ - \pi_p(C_M u_\varepsilon^+) \|_{L^\infty(\Omega)} \leq Ce^{-\sigma p}$$

for appropriate $\sigma > 0$. Finally, for the remainder, Theorem 1.1 yields with our choice of μ

$$\| r_M \|_{L^\infty(\Omega)} \leq C\varepsilon(\varepsilon\mu\kappa p K)^M \leq C\varepsilon(\varepsilon M K)^M \leq C\beta^{\mu\kappa p} \quad (3.32)$$

with $\beta < 1$. Thus, the remainder r_M is exponentially small on Ω , and an appropriate application of Lemma 3.5 allows us to conclude the proof of Theorem 3.3 by observing that

$$\| r_M - \pi_p(r_M) \|_{L^\infty(\Omega)} \leq C\beta^{\mu\kappa p} = Ce^{-\sigma p}, \quad \sigma = \mu\kappa |\ln \beta|.$$

□

3.3.2 Approximation of shocks

In this section, we want to demonstrate that the ideas of the “two-element” mesh of the preceding section can be used for the approximation of solutions u_ε which are smeared-out shocks. To that end, let us consider

$$L_\varepsilon u_\varepsilon = f + \delta_0 \quad \text{on } \Omega, \quad u_\varepsilon(\pm 1) = \alpha^\pm \quad (3.33)$$

where δ_0 denotes the Dirac distribution concentrated at the point $x = 0$. The following analog of Theorem 3.3 holds.

Theorem 3.7 *Let u_ε be the solution of (3.33) and assume that the coefficients a , b and the right hand side f satisfy (1.3), (1.5)–(1.7). For every ε , $\kappa > 0$ let the degree vector \vec{p} and the “four-element” mesh $\mathcal{T} = \mathcal{T}_{\kappa,\varepsilon}^4$ be given by*

$$\begin{aligned} \vec{p} &= \{p, p, p, p\}, & \mathcal{T}_{\kappa,\varepsilon}^4 &= \{I_1, I_2, I_3, I_4\} & \text{if } \kappa p \varepsilon < 1/2, \\ \vec{p} &= \{p\}, & \mathcal{T}_{\kappa,\varepsilon}^4 &= \{\Omega\} & \text{if } \kappa p \varepsilon \geq 1/2. \end{aligned} \quad (3.34)$$

where

$$I_1 = (-1, -\kappa p \varepsilon), \quad I_2 = (-\kappa p \varepsilon, 0), \quad I_3 = (0, 1 - \kappa p \varepsilon), \quad I_4 = (1 - \kappa p \varepsilon, 1).$$

Then there is $\varepsilon_0 > 0$ depending only on the constants (1.3), (1.5)–(1.7) such that for every $0 < \varepsilon \leq \varepsilon_0$ the following holds. Then there is a constant κ_0 also depending only on the constants of (1.3), (1.5)–(1.7) such that for every $0 < \kappa < \kappa_0$ there are $C, \sigma > 0$ independent of p and ε such that

$$\inf_{v_p \in S_0^{\vec{p},1}(\mathcal{T}_{\kappa,\varepsilon}^4)} \|u_\varepsilon - v_p\|_{L^\infty(\Omega)} \leq C e^{-\sigma p} \quad \forall p \in \mathbb{N}.$$

Proof: Let us first define a function u_δ with the property that $L_\varepsilon u_\delta = \delta_0$. To that end, let us introduce the following two auxiliary functions, u_L, u_R :

$$\begin{aligned} L_\varepsilon u_L &= 0 & \text{on } (-1, 0), & \quad u_L(-1) = 0, \quad u_L(0) = 1 \\ u_R &= e^{-\Lambda(x) + \Lambda(0)} & \text{on } (0, 1) \end{aligned}$$

where the function Λ is defined in (1.8). Note that $u_L(0) = u_R(0)$ and that u_R is smooth and independent of ε . In particular, $u'_R(0) = -\lambda(0)$ independently of ε . Lemma B.3 gives the existence of $C_1, C_2 > 0$ independent of ε such that $C_1 \varepsilon^{-1} \leq u'_L(0) \leq C_2 \varepsilon^{-1}$. Defining

$$D_\varepsilon := -(\varepsilon u'_R(0) - \varepsilon u'_L(0)) \quad (3.35)$$

we see that there are $C'_1, C'_2 > 0$ independent of ε such that

$$0 < C'_1 \leq D_\varepsilon \leq C'_2 \quad \forall 0 < \varepsilon \leq \varepsilon_0 \quad (3.36)$$

provided that ε_0 sufficiently small. Define now

$$u_\delta := -\frac{1}{D_\varepsilon} \begin{cases} u_L(x) & \text{if } -1 \leq x \leq 0 \\ u_R(x) & \text{if } 0 \leq x \leq 1 \end{cases} \quad (3.37)$$

Then u_δ is continuous on $\overline{\Omega}$, satisfies $L_\varepsilon u_\delta = 0$ on $(-1, 0) \cup (0, 1)$, and the jump of $-\varepsilon u'_\delta$ at $x = 0$ is 1 by the choice of D_ε . Thus, u_δ satisfies $L_\varepsilon u_\delta = \delta_0$ in the sense of distributions. By superposition, the solution u_ε of (3.33) can be written as

$$u_\varepsilon = u_\delta + \tilde{u}_\varepsilon$$

where \tilde{u}_ε solves the auxiliary problem

$$L_\varepsilon \tilde{u}_\varepsilon = f \quad \text{on } \Omega, \quad \tilde{u}_\varepsilon(-1) = \alpha^-, \quad \tilde{u}_\varepsilon(1) = \alpha^+ - u_\delta(1) \quad (3.38)$$

By Theorem 3.3, the function u_L can be approximated with the desired exponential accuracy on a two-element mesh on $(-1, 0)$. Such a two-element mesh is contained in the “four-element” mesh considered here. The function u_R is analytic and independent of ε and thus can be approximated well by polynomials on $(0, 1)$. Noting that the factor $1/D_\varepsilon$ appearing in the definition of u_δ can be bounded uniformly in ε (for $\varepsilon \leq \varepsilon_0$) finishes the approximation argument for u_δ . By this uniform bound on $1/D_\varepsilon$ and by the independence of u_R of ε , we conclude that $u_\delta(1)$ can be bounded uniformly in ε and thus Theorem 3.3 allows us to approximate \tilde{u}_ε to the desired accuracy on a two-element mesh for Ω . Such a two-element mesh is contained in the “four-element mesh” considered here which concludes the proof. □

4 Approximate test functions

Theorem 3.1 shows that the use of the upwinded test space $S_L^{\vec{p}}$ in (3.2) gives rise to a stable numerical scheme. Unfortunately, however, the shape functions $\psi_{k,j}$ in (3.8), (3.9) are themselves solutions of (local) convection-diffusion problems. For the case $p = 1$ and constant coefficients a, b , these upwinded test functions can be computed explicitly and lead to the so-called “Hemker test functions” [7]. For non-constant coefficients, however, they are not explicitly available. We show therefore now that stability can be retained even if the $\psi_{k,j}$ are known only approximately. The perturbation analysis of Section 4.1 shows that fairly weak accuracy requirements on the test functions ψ_{ij} suffice to ensure stability of the FEM. Especially for low p rather “crude” approximations to the L -splines ψ_{ij} are sufficient; this is the reason why techniques such as α -quadratic upwinding ([2]; see also [12] for an up-to-date account on these methods) and the use of Hemker test functions ([5, 6]) obtained by freezing coefficients lead to stable FEM. All these methods are in fact covered by our perturbation analysis.

4.1 Stability with approximate test functions

We introduce the approximate test space

$$\tilde{S}_L^{\vec{p}} = \text{span} \left\{ \tilde{\psi}_{k,j} : k = -1, j = 2, \dots, N \text{ and } k = 0, \dots, p_j - 2, j = 1, \dots, N \right\} \quad (4.1)$$

where the approximate test functions $\tilde{\psi}_{k,j} \in H_0^1(\Omega) \cap H^2(I_j)$, $j = 1, \dots, N$ are assumed to satisfy:

$$\begin{aligned} L_\varepsilon^* \tilde{\psi}_{-1,j} &= \eta_{-1,j} \text{ in } I_{j-1} \cup I_j, j = 2, \dots, N, \\ \tilde{\psi}_{-1,j} &= 0 \text{ elsewhere,} \\ \tilde{\psi}_{-1,j}(x_k) &= \delta_{j,k+1}, \quad k = 1, \dots, N, \end{aligned} \quad (4.2)$$

and

$$\begin{aligned} L_\varepsilon^* \tilde{\psi}_{i,j} &= L_i \left(2 \frac{x - m_j}{h_j} \right) + \eta_{i,j} \text{ in } I_j, i = 0, \dots, p_j - 2, j = 1, \dots, N, \\ \tilde{\psi}_{i,j} &= 0 \text{ elsewhere.} \end{aligned} \quad (4.3)$$

The question how to obtain such approximate test functions will be addressed below.

We show first that the bilinear form $B_{\mathcal{T}}(\cdot, \cdot)$ is stable on $S_0^{\vec{p}} \times \tilde{S}_L^{\vec{p}}$ provided the residuals $\eta_{i,j}$ in (4.2), (4.3) are sufficiently small.

Theorem 4.1 *Assume that ρ_j is as in (3.5) and that the approximate test functions $\tilde{\psi}_{k,j}$ satisfy (4.2), (4.3) with $\eta_{k,j}$ such that*

$$\Lambda_1 \leq c \min \left\{ \frac{1}{4C_1}, \gamma_M^{-4} \right\}, \quad \Lambda_2 \leq c \min \left\{ \frac{1}{10}, \gamma_M^{-4} \right\} \quad (4.4)$$

where $c < 1$,

$$\Lambda_1 := \sum_{j=1}^N \|\eta_{-1,j}\|_{L^2(I_j)}^2 + \|\eta_{-1,j+1}\|_{L^2(I_j)}^2, \quad (4.5)$$

$$\Lambda_2 := \max_{1 \leq j \leq N} \left\{ h_j^{-1} \sum_{i=0}^{p_j-2} (2i+1) \|\eta_{ij}\|_{L^2(I_j)}^2 \right\} \quad (4.6)$$

(we set $\Lambda_2 := 0$ if $p = 1$ and $\eta_{-1,1} = \eta_{-1,N+1} = 0$) and C_1 is the constant in (2.12).

There exists $C > 0$ independent of ε , \vec{p} and \mathcal{T} such that

$$\inf_{0 \neq v \in \tilde{S}_L^{\vec{p}}} \sup_{0 \neq u \in S^{\vec{p}}} \frac{B_{\mathcal{T}}(u, v)}{\|u\|_{H_0^0} \|v\|_{H_0^2}} \geq \frac{C}{\gamma_M} > 0. \quad (4.7)$$

Proof: Let $\tilde{v} \in \tilde{S}_L^{\vec{p}}$ be given. Then

$$\tilde{v}(x) = \sum_{j=2}^N \tilde{v}(x_{j-1}) \tilde{\psi}_{-1,j}(x) + \sum_{j=1}^N \sum_{i=0}^{p_j-2} b_{ij} \tilde{\psi}_{ij}(x).$$

We select $u_{\tilde{v}}$ as in the proof of Theorem 3.1, i.e.

$$u_{\tilde{v}}|_{I_j} = \sum_{i=0}^{p_j} a_{ij} L_i \left(2 \frac{x - m_j}{h_j} \right)$$

where

$$a_{ij} = b_{ij} \quad i = 0, \dots, p_j - 2$$

and $a_{p_j-1,j}, a_{p_j,j}$ are selected as in (3.15)-(3.17), with \tilde{v} in place of v . Then $u_{\tilde{v}}$ is continuous on $[-1, 1]$. With the test functions ψ_{ij} in (3.8), (3.9) we define also

$$v(x) := \sum_{j=2}^N \tilde{v}(x_{j-1})\psi_{-1,j}(x) + \sum_{j=1}^N \sum_{i=0}^{p_j-2} b_{ij}\psi_{ij}(x)$$

and we set

$$\delta v := \tilde{v} - v = \sum_{j=2}^N \tilde{v}(x_{j-1})\eta_{-1,j}(x) + \sum_{j=1}^N \sum_{i=0}^{p_j-2} b_{ij}\eta_{ij}(x).$$

Then

$$\begin{aligned} B_{\mathcal{T}}(u_{\tilde{v}}, \tilde{v}) &= \sum_{j=1}^N \int_{I_j} u_{\tilde{v}} L_{\varepsilon}^* v dx - \sum_{j=1}^{N-1} u_{\tilde{v}}(x_j) [\varepsilon \tilde{v}'(x_j)] \\ &+ \sum_{j=1}^N \int_{I_j} u_{\tilde{v}} L_{\varepsilon}^* \delta v dx. \end{aligned}$$

We calculate

$$\int_{I_j} u_{\tilde{v}} L_{\varepsilon}^* v dx = h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} = h_j S_j$$

and, by (3.15), (3.16),

$$- \sum_{j=1}^{N-1} u_{\tilde{v}}(x_j) [\varepsilon \tilde{v}'(x_j)] = \sum_{j=1}^{N-1} \rho_j^{-1} |[\varepsilon \tilde{v}'(x_j)]|^2,$$

hence we find

$$B_{\mathcal{T}}(u_{\tilde{v}}, \tilde{v}) \geq \sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \rho_j^{-1} |[\varepsilon \tilde{v}'(x_j)]|^2 - \sum_{j=1}^N \|u_{\tilde{v}}\|_{L^2(I_j)} \|L_{\varepsilon}^* \delta v\|_{L^2(I_j)}. \quad (4.8)$$

Now

$$\|u_{\tilde{v}}\|_{L^2(I_j)}^2 = h_j \sum_{i=0}^{p_j-2} \frac{|b_{ij}|^2}{2i+1} + h_j \sum_{i=p_j-1}^{p_j} \frac{|a_{ij}|^2}{2i+1}.$$

Reasoning as in the proof of Theorem 3.1, we find then

$$\sum_{j=1}^N \|u_{\tilde{v}}\|_{L^2(I_j)}^2 \leq 4 \sum_{j=1}^{N-1} \frac{|[\varepsilon \tilde{v}'(x_j)]|^2}{\rho_j} + (p+3) \sum_{j=1}^N h_j S_j. \quad (4.9)$$

Consider now $\|L_{\varepsilon}^* \delta v\|_{L^2(I_j)}$. We have by (4.2), (4.3)

$$(L_{\varepsilon}^* \delta v)|_{I_j} = \tilde{v}(x_{j-1})\eta_{-1,j} + \tilde{v}(x_j)\eta_{-1,j+1} + \sum_{i=0}^{p_j-2} b_{ij}\eta_{ij}.$$

Using (2.12), we estimate

$$\begin{aligned} \|L_{\varepsilon}^* \delta v\|_{L^2(I_j)} &\leq \|\tilde{v}\|_{L^{\infty}} \left(\|\eta_{-1,j}\|_{L^2(I_j)} + \|\eta_{-1,j+1}\|_{L^2(I_j)} \right) + \sum_{i=0}^{p_j-2} b_{ij} \|\eta_{ij}\|_{L^2(I_j)} \\ &\leq \|\tilde{v}\|_{L^{\infty}} \left(\|\eta_{-1,j}\|_{L^2(I_j)} + \|\eta_{-1,j+1}\|_{L^2(I_j)} \right) \\ &+ (h_j S_j)^{1/2} \left(h_j^{-1} \sum_{i=0}^{p_j-2} (2i+1) \|\eta_{ij}\|_{L^2(I_j)}^2 \right)^{1/2} \end{aligned}$$

i.e.

$$\|L_\varepsilon^* \delta v\|_{L^2(I_j)}^2 \leq 4 \|\tilde{v}\|_{L^\infty}^2 \Lambda_{1j} + h_j S_j \Lambda_{2j} \quad j = 1, \dots, N, \quad (4.10)$$

where we defined

$$\Lambda_{1j} := \|\eta_{-1,j}\|_{L^2(I_j)}^2 + \|\eta_{-1,j+1}\|_{L^2(I_j)}^2$$

and

$$\Lambda_{2j} := h_j^{-1} \sum_{i=0}^{p_j-2} (2i+1) \|\eta_{ij}\|_{L^2(I_j)}^2.$$

Hence we may estimate with (4.9)

$$\begin{aligned} & \sum_{j=1}^N \|u_{\tilde{v}}\|_{L^2(I_j)} \|L_\varepsilon^* \delta v\|_{L^2(I_j)} \\ & \leq \left(4 \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} + (p+3) \sum_{j=1}^N h_j S_j \right)^{1/2} \left(4\Lambda_1 \|\tilde{v}\|_{L^\infty}^2 + \Lambda_2 \sum_{j=1}^N h_j S_j \right)^{1/2} \\ & \leq \max\{4, p+3\}^{1/2} \left(\sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \right)^{1/2} \left(4\Lambda_1 \|\tilde{v}\|_{L^\infty}^2 + \Lambda_2 \sum_{j=1}^N h_j S_j \right)^{1/2} \end{aligned} \quad (4.11)$$

With (4.10), the embedding (2.12) and the definition of the $H_{\mathcal{T}}^2$ norm we get further

$$\begin{aligned} \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 & \leq 2 \sum_{j=1}^N \|L_\varepsilon^* v\|_{L^2(I_j)}^2 + 2 \|L_\varepsilon^* \delta v\|_{L^2(I_j)}^2 + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \\ & = 2(1 + \Lambda_2) \sum_{j=1}^N h_j S_j + 8C_1 \Lambda_1 \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \end{aligned}$$

and, after regrouping terms, it follows that

$$\sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \geq D(\eta) \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 \quad (4.12)$$

provided Λ_1, Λ_2 are sufficiently small and

$$D(\eta) := \frac{1 - 8C_1 \Lambda_1}{2(1 + \Lambda_2)}.$$

The inequality which is converse to (4.12) also holds. To obtain it, we proceed as follows: we estimate

$$\begin{aligned} \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 & = \sum_{j=1}^N \|L_\varepsilon^* v + L_\varepsilon^* \delta v\|_{L^2(I_j)}^2 + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \\ & \geq \left(\frac{1}{2} - 5\Lambda_2 \right) \sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} - 20C_1 \Lambda_1 \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 \end{aligned}$$

and obtain after rearranging terms

$$\sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \leq C(\eta) \|\tilde{v}\|_{H_{\mathcal{T}}^2}^2 \quad (4.13)$$

where

$$C(\eta) := \frac{1 + 20C_1\Lambda_1}{1/2 - 5\Lambda_2}.$$

From (4.13) and (4.9) we find

$$\|u_{\tilde{v}}\|_{H_T^0}^2 \leq \gamma_M^2 \left\{ \sum_{j=1}^N h_j S_j + \sum_{j=1}^{N-1} \frac{|\varepsilon \tilde{v}'(x_j)|^2}{\rho_j} \right\} \leq \max\{5, p+3\} C(\eta) \|\tilde{v}\|_{H_T^2}^2. \quad (4.14)$$

Further, (4.8) and (4.11) imply with (4.12) and (4.13) the bound

$$B_{\mathcal{T}}(u_{\tilde{v}}, \tilde{v}) \geq \frac{1}{\gamma_M} \left[\frac{D(\eta)}{(C(\eta))^{1/2}} - \gamma_M^2 \left(4C_1 \frac{D(\eta)}{C(\eta)} \Lambda_1 + D(\eta) \Lambda_2 \right)^{1/2} \right] \|u_{\tilde{v}}\|_{H_T^0} \|\tilde{v}\|_{H_T^2}$$

from where the inf-sup condition (4.7) follows. □

Remark 4.2 The test functions $\tilde{\psi}_{ij}$ in (4.2),(4.3) are *conforming*, i.e., globally in $H_0^1(\Omega)$ and elementwise in $H^2(I_j)$. As we shall show shortly, it is possible to obtain numerical approximations $\tilde{\psi}_{ij}$ by solving the problems (3.8), (3.9) with a least squares FEM on a subgrid $\tilde{\mathcal{T}}$ of \mathcal{T} . The assumption $\tilde{\psi}_{ij} \in H^2(I_j)$ for $I_j \in \mathcal{T}$ then implies that the least squares FEM must be *locally* C^1 conforming. Although this can be achieved, we can also admit C^0 -approximations $\tilde{\psi}_{ij}$ in Theorem 4.1, if we penalize the flux-jumps of $\tilde{\psi}_{ij}$ on the subgrid appropriately. This will complicate the following analysis, but does not pose any essential difficulties.

The stability (4.7) together with the fact that the perturbed test functions are H_T^2 -conforming and with Propositions A.4, A.3 and the approximation property Theorem 3.3 imply the following convergence result.

Theorem 4.3 *For any mesh \mathcal{T} the hp FE-solution $\tilde{u}_M \in S_0^{\vec{p}}(\mathcal{T})$ in (3.4) corresponding to the approximate test space $\tilde{S}_L^{\vec{p}}$ defined in (4.1) - (4.3) and satisfying (4.4), exists and is quasi-optimal, i.e., with C, γ_M of (4.7)*

$$\|u - \tilde{u}_M\|_{H_T^0} \leq \left(1 + \frac{\gamma_M}{C}\right) \|u - v\|_{H_T^0} \quad \forall v \in S_0^{\vec{p}}. \quad (4.15)$$

In particular, if the coefficients a, b , and the right hand side f are analytic and satisfy (1.3), (1.5)–(1.7) and the mesh $\mathcal{T} = \mathcal{T}_{\kappa, \varepsilon}$ is chosen as in (3.24) with κ sufficiently small, we have robust exponential convergence, i.e.,

$$\|u - \tilde{u}_M\|_{H_T^0} \leq C \exp(-\theta M) \quad (4.16)$$

where $C, \theta > 0$ are independent of ε, p .

Remark 4.4 The meshes $\mathcal{T}_{\kappa, \varepsilon}$ considered in Theorem 3.3 are essentially the “minimal” meshes that can resolve the boundary layer behavior of the solution u_ε in a p -version setting at a robust exponential rate. Clearly, the approximation results of the form (3.25)

hold true for any mesh \mathcal{T} that contains one small element of size $O(\varepsilon p)$ at the outflow boundary, i.e., if $\mathcal{T}_{\kappa,\varepsilon} \subset \mathcal{T}$. Due to the quasi-optimality (4.15), error estimates analogous to (4.16) hold for all meshes \mathcal{T} with $\mathcal{T}_{\kappa,\varepsilon} \subset \mathcal{T}$. These “minimal” meshes $\mathcal{T}_{\kappa,\varepsilon}$ depend on the polynomial degree p . In practice, it may be more convenient to fix a mesh \mathcal{T} and then increase p until the desired accuracy is reached. For example, piecewise polynomials on a mesh which is graded geometrically towards the outflow boundary have approximation properties similar to the minimal meshes $\mathcal{T}_{\kappa,\varepsilon}$ provided that the small element of the geometric mesh is $O(\varepsilon)$ (cf. [9]).

4.2 Computation of the Approximate Test Functions

To obtain approximate test functions $\tilde{\psi}_{ij}$, many strategies are possible: classical approaches use approximate analytical expressions (e.g., the Hemker test functions or the α -quadratic upwinding mentioned above) or asymptotic expansions (e.g., in [19]). These semianalytical approaches work well for low order methods and one-dimensional problems; for two-dimensional problems, however, and to accommodate high p together with arbitrary meshes, a fully numerical method for the computation of the test functions seems to be desirable.

Here we propose and analyze a *local least squares FEM* to approximate the ψ_{ij} stably and completely computationally. The approach allows moreover for controlling the quantities Λ_1, Λ_2 in (4.5), (4.6) a-posteriori.

The plan for the remainder of this section is as follows. In Section 4.2.1 we will define the local least squares problems which define approximate test functions $\tilde{\psi}_{-1,j}, \tilde{\psi}_{ij}$ by minimizing appropriate quadratic functionals over finite dimensional spaces \mathcal{A}_j^{hq} . In this framework, the choice of the spaces \mathcal{A}_j^{hq} determines completely the test functions $\tilde{\psi}_{-1,j}, \tilde{\psi}_{ij}$ and thus the method (3.4). The exact test functions are also solutions of singularly perturbed convection-diffusion equations with analytic coefficients. We will therefore choose \mathcal{A}_j^{hq} as spaces of piecewise polynomials of degree q on a two-element mesh (one small element at the outflow boundary of the local problem and one large element) in complete analogy to our approximation theory for the global solution u_ε . The details of these approximation results are provided in Sections 4.2.2, 4.2.3.

Other choices of the spaces \mathcal{A}_j^{hq} lead to different methods. For example, for $p = 1$ and \mathcal{A}_j^{hq} consisting of quadratic polynomials, the least squares method yields approximate test functions $\tilde{\psi}_{-1,j}$ very similar to those obtained by α -quadratic upwinding. The Hemker test functions for $p = 1$ and constant coefficients a, b may be obtained with our least squares method if one includes in the spaces \mathcal{A}_j^{hq} exponentials which solve the homogeneous adjoint problem.

4.2.1 Approximate Test Functions via Local Least Squares FEM

To motivate the method, we define $\mathcal{A}_j := (H^2 \cap H_0^1)(I_j)$, $j = 1, \dots, N$. We define further $\varphi_j(x)$ to be the piecewise linear “hat” function with

$$\varphi_j(x_{j-1}) = \varphi_j(x_{j+1}) = 0, \varphi_j(x_j) = 1, \varphi_j(x) = 0 \text{ on } \Omega \setminus \overline{I_{j-1} \cup I_j}.$$

Then $\psi_{-1,j} - \varphi_{j-1} \in \mathcal{A}_{j-1} \cup \mathcal{A}_j$ and we have the variational characterization

$$(\psi_{-1,j} - \varphi_{j-1})|_{I_k} = \arg \min_{\psi \in \mathcal{A}_k} \|L_\varepsilon^*(\psi - \varphi_{j-1})\|_{L^2(I_k)}^2, k = j-1, j, \quad (4.17)$$

for $j = 2, \dots, N$ and

$$\psi_{ij}|_{I_j} = \arg \min_{\psi \in \mathcal{A}_j} \left\| L_\varepsilon^* \psi - L_i \left(2 \frac{\cdot - m_j}{h_j} \right) \right\|_{L^2(I_j)}^2 \quad i = 0, \dots, p_j - 2, j = 1, \dots, N. \quad (4.18)$$

The assumptions (1.3), (1.4) imply that the operator L_ε and therefore also its adjoint L_ε^* are injective from $\mathcal{A}_j \rightarrow L^2(I_j)$. Hence the expression

$$\|\psi\|_{*,j} := \|L_\varepsilon^* \psi\|_{L^2(I_j)} \quad (4.19)$$

is a norm on \mathcal{A}_j (homogeneity and triangle inequality being obvious) and the quadratic functionals in (4.17), (4.18) are strictly convex and lower semicontinuous. Therefore (4.17), (4.18) admit unique solutions $\psi_{-1,j}$, ψ_{ij} which coincide with those in (3.8), (3.9).

For a numerical approximation of the functions $\psi_{-1,j}$, ψ_{ij} , let $\mathcal{A}_j^{hq} \subset \mathcal{A}_j$ be a finite dimensional subspace. We obtain external approximate test functions $\tilde{\psi}_{-1,j}$ by

$$(\tilde{\psi}_{-1,j} - \varphi_{j-1})|_{I_k} = \arg \min_{\psi \in \mathcal{A}_k^{hq}} \|L_\varepsilon^*(\psi - \varphi_{j-1})\|_{L^2(I_k)}^2, k = j-1, j \quad (4.20)$$

for $j = 2, \dots, N$ and internal approximate test functions $\tilde{\psi}_{ij}$, $i = 0, \dots, p_j - 2$ by

$$\tilde{\psi}_{ij}|_{I_j} = \arg \min_{\psi \in \mathcal{A}_j^{hq}} \left\| L_\varepsilon^* \psi - L_i \left(2 \frac{\cdot - m_j}{h_j} \right) \right\|_{L^2(I_j)}^2 \quad j = 1, \dots, N. \quad (4.21)$$

These approximations are also uniquely defined. Moreover, they are optimal in the norm $\|\cdot\|_{*,j}$, for we have

$$\|\psi_{-1,j} - \tilde{\psi}_{-1,j}\|_{*,k} \leq \|\psi_{-1,j} - \varphi_{j-1} - \psi\|_{*,k} \quad \forall \psi \in \mathcal{A}_k^{hq}, k = j-1, j, \quad j = 2, \dots, N \quad (4.22)$$

$$\|\psi_{ij} - \tilde{\psi}_{ij}\|_{*,j} \leq \|\psi_{ij} - \psi\|_{*,j}, \quad \forall \psi \in \mathcal{A}_j^{hq}, \quad j = 1, \dots, N. \quad (4.23)$$

Thus, the design of the approximation spaces \mathcal{A}_j^{hq} proceeds in the usual fashion: based on the regularity of the exact test functions ψ_{ij} , we show that we can select the least squares approximation spaces \mathcal{A}_j^{hq} so that exponential convergence rates in the global $\|\psi\|_* := \|L_\varepsilon^* \psi\|_{L^2(\Omega)}$ norm can be achieved.

Remark 4.5 The calculation of the approximate test functions $\tilde{\psi}_{-1,j}$, $\tilde{\psi}_{ij}$ can be done efficiently if one observes that eqs. (4.20), (4.21) represent completely decoupled local problems on the elements I_j , which can be solved in parallel. Furthermore, on each element I_j the local least squares problems (4.20), (4.21) can be solved efficiently because the equivalent matrix formulations lead to problems with the same stiffness matrix and merely different right hand sides. Thus, once a convenient decomposition of the elemental least squares matrix is found (e.g., its LU decomposition), the approximate test functions $\tilde{\psi}_{-1,j}|_{I_j}$, $\tilde{\psi}_{-1,j+1}|_{I_j}$, $\tilde{\psi}_{ij}$, $i = 0, \dots, p_j - 2$ can be obtained by $p_j + 1$ backsolves.

4.2.2 Approximation Results on the Reference Element

We begin our analysis with the following approximation result, quite analogous to Theorem 3.3. Introduce the ε dependent norm $\|\cdot\|_{2,\varepsilon}$ on $H^2(I)$ by

$$\|u\|_{2,\varepsilon} := \varepsilon \|u''\|_{L^2(I)} + \|u'\|_{L^2(I)} + \|u\|_{L^2(I)} \quad (4.24)$$

Theorem 4.6 *Let u_ε be the solution of (1.1), (1.2). Assume that (1.3), (1.5)–(1.7) hold and let \vec{p} and $\mathcal{T}_{\kappa,\varepsilon}$ be as in Theorem 3.3. Then there is $\kappa_0 > 0$ depending only on the constants of (1.3), (1.5)–(1.7) such that for every $0 < \kappa < \kappa_0$ there are $C, \sigma > 0$ independent of p, ε and functions $v_p \in S^{\vec{p},2}(\mathcal{T}_{\kappa,\varepsilon})$ such that*

$$u_\varepsilon(\pm 1) = v_p(\pm 1), \quad \|u_\varepsilon - v_p\|_{2,\varepsilon} \leq C\varepsilon^{-1/2}e^{-\sigma p} \quad \forall p \in \mathbb{N}. \quad (4.25)$$

Proof: The proof is very similar to the proof of Theorem 3.3. We will therefore only highlight the main differences. For the asymptotic case $\kappa p \varepsilon \geq 1$, the claim follows easily as in the proof of Theorem 3.3. Let us consider the pre-asymptotic case $\kappa p \varepsilon < 1$. Choose μ, M as in (3.28) and use the decomposition (1.9) for u_ε . We get from Theorem 1.1 and Lemma 3.6

$$\|(w_M - i_p w_M)^{(l)}\|_{L^\infty(I)} \leq C e^{-\sigma p}, \quad l = 0, 1, 2.$$

For the remainder r_M , we use Theorem 1.1 to get for $0 \leq n \leq 2$ and the implicit assumption on p that $M \geq 1$:

$$\|r_M^{(n)}\|_{L^\infty(I)} \leq C\varepsilon^{1-n}(\varepsilon\mu\kappa p M)^M \leq C\varepsilon^{1-n}(\varepsilon\mu\kappa p)^{n-1}(\varepsilon\mu\kappa p M)^{M-n+1} \leq C(\mu\kappa p)^{n-1}q^{M-n+1}.$$

Thus, an application of Lemma 3.5 yields

$$\|(r_M - i_p r_M)^{(l)}\|_{L^\infty(I)} \leq e^{-\sigma p}, \quad l = 0, 1, 2$$

for appropriately chosen $\sigma > 0$. Let us now turn to the approximation of the boundary layer function u_ε^+ . The technical details are very similar to the proof of Theorem 5.1 of [15]. As in the proof of Theorem 3.3, let l_1, l_2 be the two linear maps from the reference element I onto the two elements I_1, I_2 . Consider first the approximation of the function

$$U(x) := (u_\varepsilon^+)'$$

and let $\tilde{x} := 1 - \kappa p \varepsilon$ be the internal node of the mesh $\mathcal{T}_{\kappa,\varepsilon}$. Note that, up to a factor ε^{-1} , the function U satisfies similar estimates as u_ε^+ by Theorem 1.1. Define the approximant $U_p \in S^{\vec{p},1}(\mathcal{T}_{\kappa,\varepsilon})$ by

$$U_p(x) := \begin{cases} U(-1) + \frac{\varepsilon^{1/2}U(\tilde{x}) - U(-1)}{\tilde{x}+1}(x+1) & \text{on } I_1 \\ (i_p(U \circ l_2)) \circ l_2^{-1} - \frac{(1-\varepsilon^{1/2})U(\tilde{x})}{\kappa p \varepsilon}(1-x) & \text{on } I_2. \end{cases}$$

We claim now that for κ sufficiently small, there are $C, \sigma > 0$ such that

$$\|(U - U_p)^{(l)}\|_{L^2(I)} \leq C\varepsilon^{-1/2-l}e^{-\sigma p}, \quad l = 0, 1, \quad \forall p \in \mathbb{N}. \quad (4.26)$$

For the approximation on the small boundary layer element I_2 , we calculate analogously to (3.30) (after absorbing the powers of p arising from the use of Lemma 3.6 and choosing κ_0 sufficiently small)

$$\| (U \circ l_2 - (i_p(U \circ l_2)))^{(l)} \|_{L^\infty(I)} \leq C\varepsilon^{-1}e^{-\sigma p}, \quad l = 0, 1.$$

Hence,

$$\| (U - (i_p(U \circ l_2) \circ l_2^{-1}))^{(l)} \|_{L^2(I_2)} \leq C\varepsilon^{-1}(\kappa p \varepsilon)^{1/2-l}e^{-\sigma p}, \quad l = 0, 1.$$

We calculate further with $0 < \varepsilon \leq 1$:

$$\| \left(\frac{(1 - \varepsilon^{1/2})U(\tilde{x})}{\kappa p \varepsilon} (1 - x) \right)^{(l)} \|_{L^2(I_2)} \leq |U(\tilde{x})|(\kappa p \varepsilon)^{1/2-l}, \quad l = 0, 1.$$

Combining these two last estimates and observing that by Theorem 1.1 $|U(\tilde{x})| \leq C\varepsilon^{-1}e^{-\underline{a}\kappa p/2}$, we get for some suitable $\sigma > 0$

$$\| (U - U_p)^{(l)} \|_{L^2(I_2)} \leq C\varepsilon^{-1/2-l}e^{-\sigma p} \quad l = 0, 1.$$

Let us now consider the large element I_1 . We have

$$\| (U - U_p)^{(l)} \|_{L^2(I_1)} \leq \|U^{(l)}\|_{L^2(I_1)} + \|U_p^{(l)}\|_{L^2(I_1)}$$

By Theorem 1.1, $\|U^{(l)}\|_{L^2(I_1)} \leq C\varepsilon^{-1/2-l}e^{-\underline{a}\kappa p/2}$, and it is easy to verify that

$$\|U_p^{(l)}\|_{L^2(I_1)} \leq C \max(\varepsilon^{1/2}|U(\tilde{x})|, |U(-1)|) \leq C\varepsilon^{-1/2}e^{-\underline{a}\kappa p/2}, \quad l = 0, 1.$$

Hence, we have proven (4.26). To conclude the proof of Theorem 4.6, we define the approximant $v_{p+1} \in S^{\bar{p}+1,2}(\mathcal{T}_{\kappa,\varepsilon})$ by

$$v_{p+1}(x) := u_\varepsilon(-1) + \int_{-1}^x U_p(t) dt - \frac{1}{2} \left\{ \int_{-1}^1 U_p(t) - U(t) dt \right\} (x + 1).$$

We note $v_{p+1}(\pm 1) = u_\varepsilon(\pm 1)$ and the observation

$$\begin{aligned} \left| \int_{-1}^1 U_p(t) - U(t) dt \right| &\leq C\varepsilon^{-1/2}e^{-\sigma p} \\ (v_{p+1} - u_\varepsilon)(x) &= \int_{-1}^x U_p(t) - U(t) dt - \frac{1}{2} \left\{ \int_{-1}^1 U_p(t) - U(t) dt \right\} (x + 1) \end{aligned}$$

allows us to conclude the argument. \square

This result immediately applies also to FE-approximations of the adjoint problem which will then enable us to analyze the approximation of the test functions.

Corollary 4.7 *Assume (1.4), (1.5)–(1.7). For $q \in \mathbb{N}$, $\kappa > 0$ set*

$$\begin{aligned} \vec{q} = \{q, q\}, \quad \mathcal{T}_{\kappa,\varepsilon}^* = \{I_1, I_2\} \quad I_1 = [-1, -1 + \kappa q \varepsilon], I_2 = [-1 + \kappa q \varepsilon, 1] &\quad \text{if } \kappa q \varepsilon < 1, \\ \vec{q} = \{q\}, \quad \mathcal{T}_{\kappa,\varepsilon}^* = \{[-1, 1]\} &\quad \text{if } \kappa q \varepsilon \geq 1. \end{aligned} \tag{4.27}$$

Let u_ε^* be the solution of the adjoint problem

$$L_\varepsilon^* u_\varepsilon^* = f \quad \text{on } \Omega, \quad u_\varepsilon^*(\pm 1) = \alpha^\pm. \quad (4.28)$$

Then there is $\kappa_0 > 0$ depending only on the constants in (1.4), (1.5)–(1.7) and α^\pm such that the following holds. For each $0 < \kappa < \kappa_0$ there are $C, \sigma > 0$ independent of q, ε and a sequence (v_q) of functions in $S^{\bar{q},2}(\mathcal{T}_{\kappa,\varepsilon}^*)$ such that

$$v_q(\pm 1) = u_\varepsilon^*(\pm 1) \quad \text{and} \quad \|u_\varepsilon^* - v_q\|_{2,\varepsilon} \leq C\varepsilon^{-1/2}e^{-\sigma q} \quad \forall q \in \mathbb{N}$$

Proof: The change of variables $x \mapsto -x$ changes L_ε^* into a differential operator of the type of L_ε whose coefficients satisfy by (1.4) all the necessary conditions for Theorem 4.6 to imply the result. □

We will use Corollary 4.7 to estimate the errors

$$\|\psi_{ij} - \tilde{\psi}_{ij}\|_{*,j} = \|L_\varepsilon^*(\psi_{ij} - \tilde{\psi}_{ij})\|_{L^2(\Omega)}$$

of the $\tilde{\psi}_{ij}$ computed by the least squares methods (4.18), (4.21). Evidently, this will imply also bounds on Λ_1, Λ_2 in Theorem 4.1.

Clearly, Corollary 4.7 could be applied on Ω_j after a scaling argument; however, rather than the general right hand side f in (4.28), we must also solve problem (4.3) with $f = L_i$, $i = 0, \dots, p_j - 2$. To that end, let us formulate the following

Proposition 4.8 *Let $\mathcal{T}_{\kappa,\varepsilon}^*$ be as in Corollary 4.7 and let u_ε^* be the solution of (4.28) where the right hand side f is a polynomial of degree p . Then there is $\kappa_0 > 0$ depending only on the constants in (1.4), (1.5), (1.6) such that the following holds. For each $0 < \kappa < \kappa_0$ there are constants $C, \sigma, \tau > 0$ and a sequence (v_q) of functions in $S^{\bar{q},2}(\mathcal{T}_{\kappa,\varepsilon}^*)$ such that for all $q \geq \tau p$*

$$v_q(\pm 1) = u_\varepsilon(\pm 1) \quad \text{and} \quad \|u_\varepsilon - v_q\|_{2,\varepsilon} \leq C\varepsilon^{-1/2}e^{-\sigma q} \left(|\alpha^-| + |\alpha^+| + \|f\|_{L^\infty(I)} \right). \quad (4.29)$$

For the proof of Proposition 4.8 we need the following lemma.

Lemma 4.9 (Bernstein's lemma) *Let $I = [-1, 1]$. For every $\rho > 1$ there are $C_\rho, \gamma_\rho > 0$ such that for all polynomials P_p of degree p*

$$\|P_p^{(n)}\|_{L^\infty(I)} \leq C_\rho n! \gamma_\rho^n \rho^p \|P_p\|_{L^\infty(I)} \quad \forall n \in \mathbb{N}_0.$$

Proof: For $\rho > 1$ denote \mathcal{E}_ρ the ellipse (in the complex plane) whose foci are ± 1 and whose axes have lengths $\rho + \rho^{-1}, \rho - \rho^{-1}$. By Bernstein's Lemma (e.g., [8], III.15) the extension of P_p to the complex plane satisfies

$$\|P_p\|_{L^\infty(\mathcal{E}_\rho)} \leq \rho^p \|P_p\|_{L^\infty(I)} \quad \forall \rho > 1.$$

The claim of the lemma follows by Cauchy's integral theorem for derivatives. □

Proof of Proposition 4.8: By linearity, we may write the solution u_ε^* as the sum of $u_{\varepsilon,h} + u_{\varepsilon,p}$ where $u_{\varepsilon,h}$ solves (4.28) with homogeneous right hand side and inhomogeneous Dirichlet data α^\pm and where $u_{\varepsilon,p}$ solves (4.28) with right hand side f and homogeneous Dirichlet data. Corollary 4.7 implies (4.29) for $u_{\varepsilon,h}$. We may therefore concentrate on the approximation of $u_{\varepsilon,p}$. Fix $\rho > 1$. Lemma 4.9 implies

$$\|f^{(n)}\|_{L^\infty(I)} \leq (C_\rho \rho^p \|f\|_{L^\infty(I)}) \gamma_\rho^n n! \quad \forall n \in \mathbb{N}_0. \quad (4.30)$$

Consider the scaled function

$$\tilde{u}_{\varepsilon,p} := \frac{u_{\varepsilon,p}}{\rho^p \|f\|_{L^\infty(I)}}$$

which solve the equation

$$\begin{aligned} L_\varepsilon^* \tilde{u}_{\varepsilon,p} &= \tilde{f} \quad \text{on } I, \quad \tilde{u}_{\varepsilon,p}(\pm 1) = 0 \\ \|\tilde{f}^{(n)}\|_{L^\infty(I)} &\leq C_\rho \gamma_\rho^n n! \quad \forall n \in \mathbb{N}_0 \end{aligned}$$

Hence, we may apply Corollary 4.7 to the function $\tilde{u}_{\varepsilon,p}$ and obtain the existence of functions $\tilde{v}_q \in S^{\vec{q},2}(\mathcal{T}_{\kappa,\varepsilon}^*)$ with

$$\tilde{v}_q(\pm 1) = 0 \quad \|\tilde{u}_{\varepsilon,p} - \tilde{v}_q\|_{2,\varepsilon} \leq C \varepsilon^{-1/2} e^{-\sigma q} \quad \forall q \in \mathbb{N}$$

Scaling back, we obtain for the function $v_q := \rho^p \|f\|_{L^\infty(I)} \tilde{v}_q \in S^{\vec{q},2}(\mathcal{T}_{\kappa,\varepsilon}^*)$

$$v_q(\pm 1) = 0 \quad \|u_{\varepsilon,p} - v_q\|_{2,\varepsilon} \leq C \varepsilon^{-1/2} e^{-\sigma q} \rho^p \|f\|_{L^\infty(I)} \quad \forall q \in \mathbb{N}$$

As $q = q/2 + q/2 \geq q/2 + \tau p/2$ we see that choosing τ sufficiently large implies the statement of the proposition. \square

4.2.3 Analysis of the Local Least Squares FEM

In the preceding subsection, we analyzed the approximation properties of piecewise polynomials on the reference element I . In order to obtain bounds for the errors in (4.22), (4.23), let us introduce the linear transformations l_k by

$$l_k : I \rightarrow I_k, \quad \hat{x} \mapsto l_k(\hat{x}) := m_k + \frac{h_k}{2} \hat{x}, \quad k = 1, \dots, N. \quad (4.31)$$

Furthermore, for $k = 1, \dots, N$ let us set

$$\begin{aligned} \widehat{\mathcal{A}}_k &:= \mathcal{A}_k \circ l_k \\ \widehat{u} &:= u \circ l_k \quad \forall u \in \mathcal{A}_k \\ \varepsilon_k &:= \frac{2}{h_k} \varepsilon \\ \widehat{L}_{\varepsilon_k}^* &:= -\varepsilon_k \frac{\partial^2}{\partial \hat{x}^2} - a(l_k(\hat{x})) \frac{\partial}{\partial \hat{x}} + \frac{h_k}{2} [b(l_k(\hat{x})) - a'(l_k(\hat{x}))] \end{aligned} \quad (4.32)$$

Note that $\widehat{L}_{\varepsilon_k}^* \widehat{u} = \frac{h_k}{2} (L_\varepsilon^* u) \circ l_k$ for all $u \in \mathcal{A}_k$. Note also that the coefficients of $\widehat{L}_{\varepsilon_k}^*$ satisfy (1.4) (with the same γ_1^* , γ_2^*) and estimates similar to (1.5), (1.6). A straightforward calculation gives (recall (4.19))

$$\|v\|_{*,k} = \|L_\varepsilon^* v\|_{L^2(I_k)} = \sqrt{\frac{2}{h_k}} \left\| \widehat{L}_{\varepsilon_k}^* \widehat{v} \right\|_{L^2(I)} \leq C h_k^{-1/2} \|\widehat{v}\|_{2,\varepsilon_k} \quad \forall v \in \mathcal{A}_k \quad (4.33)$$

where the constant $C > 0$ depends only on the constants of (1.5), (1.6).

In order to get bounds on the expressions (4.22), (4.23), we note that for each $k = 1, \dots, N$, the functions $\widehat{\psi}_{-1,j} := \psi_{-1,j} \circ l_k$, $\widehat{\psi}_{i,k} := \psi_{i,k} \circ l_k$ satisfy

$$\begin{aligned} \widehat{L}_{\varepsilon_k}^* \widehat{\psi}_{-1,j} &= 0 & \text{on } I, & \quad \widehat{\psi}_{-1,j}(-1) = 0, \widehat{\psi}_{-1,j}(1) = 1 & \text{if } k = j-1, \quad j = 2, \dots, N, \\ \widehat{L}_{\varepsilon_k}^* \widehat{\psi}_{-1,j} &= 0 & \text{on } I, & \quad \widehat{\psi}_{-1,j}(-1) = 1, \widehat{\psi}_{-1,j}(1) = 0 & \text{if } k = j, \quad j = 2, \dots, N, \\ \widehat{L}_{\varepsilon_k}^* \widehat{\psi}_{i,k} &= \frac{h_k}{2} L_i & \text{on } I, & \quad \widehat{\psi}_{i,k}(\pm 1) = 0, & \quad i = 0, \dots, p_k - 2, \quad k = 1, \dots, N. \end{aligned}$$

Here, the functions L_i are the Legendre polynomials which satisfy $\|L_i\|_{L^\infty(I)} = 1$. With the notation of Proposition 4.8, let us choose finite dimensional subspaces $\mathcal{A}_k^{hq} \subset \mathcal{A}_k$ as

$$\mathcal{A}_k^{hq} := S_0^{\vec{q},2}(\mathcal{T}_{\kappa,\varepsilon_k}^*) \circ l_k. \quad (4.34)$$

Concerning the approximation of the functions $\widehat{\psi}_{-1,j}$, $\widehat{\psi}_{i,k}$, in the spaces $\widehat{\mathcal{A}}_k^{hq} = S_0^{\vec{q},2}(\mathcal{T}_{\kappa,\varepsilon_k}^*)$, Proposition 4.8 gives the existence of C , $\sigma > 0$, $\tau > 0$ such that

$$\begin{aligned} \inf_{\widehat{v}_q \in \widehat{\mathcal{A}}_k^{hq}} \left\| \widehat{\psi}_{-1,j} - \widehat{\varphi}_{j-1} - \widehat{v}_q \right\|_{2,\varepsilon_k} &\leq C \varepsilon_k^{-1/2} e^{-\sigma q} \quad \forall q \in \mathbb{N} \\ \inf_{\widehat{v}_q \in \widehat{\mathcal{A}}_k^{hq}} \left\| \widehat{\psi}_{i,k} - \widehat{v}_q \right\|_{2,\varepsilon_k} &\leq C \varepsilon_k^{-1/2} h_k e^{-\sigma q} \quad \forall q \geq \tau p_k \end{aligned}$$

Hence, with the choice (4.34) for the finite dimensional subspaces $\mathcal{A}_k^{hq} \subset \mathcal{A}_k$ we obtain for (4.22), (4.23) with the aid of (4.33) and $\varepsilon_k = 2\varepsilon/h_k$

$$\left\| \psi_{-1,j} - \widetilde{\psi}_{-1,j} \right\|_{*,k} \leq C \varepsilon^{-1/2} e^{-\sigma q} \quad \forall q \in \mathbb{N}, k = j-1, j, \quad j = 2, \dots, N, \quad (4.35)$$

$$\left\| \psi_{i,k} - \widetilde{\psi}_{i,k} \right\|_{*,k} \leq C h_k^{1/2} \varepsilon^{-1/2} e^{-\sigma q} \quad \forall q \geq \tau p_k, \quad k = 1, \dots, N, \quad (4.36)$$

These estimates allow us to control the expressions for Λ_1 , Λ_2 arising in Theorem 4.1. We obtain

Theorem 4.10 *Let $\mathcal{T} = \{I_1, \dots, I_N\}$ be any mesh on Ω and \vec{p} any degree vector. For $\kappa > 0$, define a subgrid mesh $\mathcal{T}_{\kappa,\varepsilon}^* := \cup_{k=1}^N \mathcal{T}_{\kappa,\varepsilon_k}^{*,k}$ where for each element I_k , the subdivision $\mathcal{T}_{\kappa,\varepsilon_k}^{*,k}$ is given by*

$$\begin{aligned} \mathcal{T}_{\kappa,\varepsilon_k}^{*,k} &= \{l_k(J_1), l_k(J_2)\} & J_1 = [-1, -1 + \kappa q \varepsilon_k], J_2 = [-1 + \kappa q \varepsilon_k, 1] & \text{if } \kappa q \varepsilon_k < 1, \\ \mathcal{T}_{\kappa,\varepsilon_k}^{*,k} &= \{I_k\} & & \text{if } \kappa q \varepsilon_k \geq 1. \end{aligned}$$

Here ε_k, h_k etc. are as in (4.32). Let furthermore the spaces \mathcal{A}_k^{hq} be given by (4.34), or, equivalently, $\mathcal{A}_k^{hq} = S_0^{\vec{q},1}(\mathcal{T}_{\kappa,\varepsilon}^*) \cap \mathcal{A}_k$ where $\vec{q} = (q, \dots, q)$.

Then, for κ sufficiently small, there exist $C, \sigma, \tau > 0$ depending only on the constants of (1.4), (1.5), (1.6) such that the approximate test functions $\tilde{\psi}_{-1,j}, \tilde{\psi}_{ij}$ of (4.20), (4.21) satisfy (4.2), (4.3) and Λ_1, Λ_2 defined in (4.5), (4.6) can be estimated as follows.

$$\begin{aligned}\Lambda_1 &\leq CN\varepsilon^{-1}e^{-\sigma q} & \forall q \in \mathbb{N}, \\ \Lambda_2 &\leq C\varepsilon^{-1}e^{-\sigma q} & \forall q \geq \tau p.\end{aligned}$$

Proof: By (4.5), (4.35) we have

$$\Lambda_1 \leq \sum_{j=2}^N \sum_{k=j-1}^j \|\psi_{-1,j} - \tilde{\psi}_{-1,j}\|_{*,k}^2 \leq CN\varepsilon^{-1}e^{-2\sigma q} \quad \forall q.$$

For Λ_2 , (4.6), (4.36) imply for all $q \geq \tau p$

$$\Lambda_2 \leq \max_{1 \leq j \leq N} \left\{ h_j^{-1} \sum_{i=0}^{p_j-2} (2i+1) \|\psi_{ij} - \tilde{\psi}_{ij}\|_{*,j}^2 \right\} \leq C \max_{1 \leq j \leq N} \varepsilon^{-1} p_j^2 e^{-2\sigma q} \leq C\varepsilon^{-1} p^2 e^{-2\sigma q}.$$

As we may assume $\tau \geq 1$, the factor p^2 can be absorbed in the exponential term at the expense of slightly reducing 2σ which concludes the proof of Theorem 4.10. \square

This result allows us finally to deduce the stability of the hp -FEM with least squares approximations of the test functions.

Corollary 4.11 *Under the hypotheses of Theorem 4.10 there is $\kappa_0 > 0$ depending only on the constants in (1.3), (1.4), (1.5)–(1.7) such that for every $0 < \kappa < \kappa_0$ there is $c > 0$ such that for $q \geq c \max(p, |\ln \varepsilon| + \ln N)$ the FEM (3.4) corresponding to \tilde{S}_L^p (computed by (4.20), (4.21)) is stable and hence quasi-optimal.*

5 Numerical Example and Implementational Aspects

The aim of the present section is to illustrate the performance of the hp FEM analyzed in this paper with particular attention to its robustness with respect to small viscosities ε . We consider the model problem

$$-\varepsilon u_\varepsilon'' + u_\varepsilon' = 1 \quad \text{on } \Omega = (-1, 1), \quad u_\varepsilon(\pm 1) = 0. \quad (5.1)$$

The exact solution has a boundary layer at the outflow boundary and is given by

$$u_\varepsilon = x + 1 + \frac{2}{1 - e^{-2/\varepsilon}} \left(e^{-2/\varepsilon} - e^{-(1-x)/\varepsilon} \right).$$

Guided by the approximation result Theorem 3.3, we choose for the trial space the spaces $S_0^{p,1}(\mathcal{T}_{\kappa,\varepsilon})$ with $\kappa = 1$ (cf. (3.1)) where the meshes $\mathcal{T}_{\kappa,\varepsilon}$ are given by (3.24). A specific basis of $S_0^{p,1}$ is given by the usual piecewise linear “nodal” shape functions and the integrated Legendre polynomials (the “internal” shape functions). For this particular problem, the basis functions $\psi_{-1,j}, \psi_{i,j}$ of the spaces of L -splines (cf. (3.8), (3.9)) can be determined in

the form of antiderivatives of exponentials times Legendre polynomials of degree up to p . Note that, since the coefficients of (5.1) are constant, the classical Hemker test functions arise for $p = 1$. Hence, for $p > 1$, in this example our scheme could be viewed as an hp version of the ‘‘Hemker test function’’ method.

Our numerical experiments were performed using MATLAB, i.e., with double precision (16 decimal) accuracy. They had the following aims: i) to show that robust exponential convergence is indeed achievable by hp -FEM, ii) to demonstrate that the method remains numerically stable as ε decreases to the order of machine precision (note that unlike e.g. the SDFEM, our method does not introduce any artificial viscosity into the computation) and iii) to assess the impact of inaccurate test functions and numerical quadrature on the stability and the robustness of the scheme.

Let us first address a few implementational aspects. If the test functions are approximated using piecewise polynomials (as proposed in Section 4), then the load vector is created in the usual way requiring some numerical integration, in general. Further, the local adjoint problems determining the (approximate) test functions may be solved completely independently and in parallel.

The stiffness matrix corresponding to $B_{\mathcal{T}}(u, v)$ consists of two parts: a mass matrix like term stemming from the domain integrals and a finite-volume like term stemming from the flux-jumps at interelement boundaries. For the mass matrix part, the test functions $\psi_{-1,j}$, $\psi_{i,j}$ need not be known completely. Rather, only $L_{\varepsilon}^* \psi_{-1,j}$, $L_{\varepsilon}^* \psi_{i,j}$ (which are chosen to be Legendre polynomials and hence known) are required. For the flux jumps, the only information employed from the Dirichlet problems (3.8), (3.9) are the normal derivatives in the endpoints, i.e., a Dirichlet-to-Neumann map is needed. In Section 4, we proposed a least squares method to approximate the test functions, but for the generation of the stiffness matrix, any sufficiently accurate Dirichlet-to-Neumann map may be taken.

For the present, constant coefficient model problem (5.1), exact representations of the test functions as antiderivatives of Legendre polynomials times a boundary layer function are available which must be integrated numerically. The next lemma shows how functions of boundary layer type can be integrated numerically in a very efficient way using standard Gaussian quadrature formulas. In our calculations, the numerical evaluation of integrals over elements of boundary layer functions against polynomials were performed based on the ideas of this ‘‘two-element’’ quadrature scheme with q points in each subelement.

Lemma 5.1 *Let w, f be analytic on $\Omega = (-1, 1)$ and satisfy*

$$\|f^{(n)}\|_{L^{\infty}(\Omega)} \leq C_f (K_f)^n n!, \quad |w^{(n)}(x)| \leq C_w (K_w)^n e^{-(1-x)/\varepsilon} \max(n, \varepsilon^{-1})^n, \quad x \in \Omega, \quad \forall n \in \mathbb{N}_0, \varepsilon \in (0, 1)$$

For $q \in \mathbb{N}$ let $\mathcal{T}_{\kappa, \varepsilon}$ be the ‘‘two-element’’ meshes introduced in (3.24) (with q taking the role of p) and denote by $G_q(\mathcal{T}_{\kappa, \varepsilon}, wf)$ the composite Gaussian quadrature rule with q points in each element applied to the function wf . Then there is $\kappa_0 > 0$ such that for $0 < \kappa < \kappa_0$ there are $C, \sigma > 0$ (independent of ε, q) such that

$$\left| \int_{\Omega} w(x)f(x) dx - G_q(\mathcal{T}_{\kappa, \varepsilon}, wf) \right| \leq C e^{-\sigma q}, \quad q = 1, 2, 3, \dots \quad (5.2)$$

If f is a polynomial of degree p with $\|f\|_{L^{\infty}(\Omega)} \leq 1$ then under the assumption $q \geq p + 1$ estimate (5.2) holds with C, σ independent of ε, p, q .

Proof: Observe that for the composite Gaussian quadrature formula of order q the quadrature error may be estimated by twice the size of the integration domain times a L^∞ best approximation of the integrand:

$$\left| \int_{\Omega} w(x)f(x) dx - G_q(\mathcal{T}_{\kappa,\varepsilon}, wf) \right| \leq 2 |\Omega| \inf_{\pi_{2q-1}} \|wf - \pi_{2q-1}\|_{L^\infty(\Omega)}$$

where the infimum is taken over all piecewise polynomials π_{2q-1} of degree $2q - 1$ on the mesh $\mathcal{T}_{\kappa,\varepsilon}$. Let $\pi_{q-1}(f)$, $\pi_q(w)$ be the piecewise Gauss-Lobatto interpolants of f , w of orders $q - 1$, q , respectively. Upon setting $\pi_{2q-1} := \pi_{q-1}(f)\pi_q(w)$ and using the stability result (3.26), we can bound

$$\|wf - \pi_{2q-1}\|_{L^\infty(\Omega)} \leq \|f - \pi_{q-1}(f)\|_{L^\infty(\Omega)}\|w\|_{L^\infty(\Omega)} + C_{GL}(1 + \ln q)\|f\|_{L^\infty(\Omega)}\|w - \pi_q(w)\|_{L^\infty(\Omega)}. \quad (5.3)$$

The proof of Theorem 3.3 shows that for the function w , which is of boundary layer type,

$$\|w - \pi_q(w)\|_{L^\infty(\Omega)} \leq Ce^{-\sigma q}$$

with C , $\sigma > 0$ independent of q , ε provided that κ is sufficiently small. A similar estimate holds for $\|f - \pi_{q-1}(f)\|_{L^\infty(\Omega)}$ by Lemma 3.6, and thus the right hand side of (5.3) is exponentially small (in q). Finally, if f is a polynomial of degree p and $q \geq p + 1$, then the term involving $f - \pi_{q-1}(f)$ vanishes in (5.3), and hence the claim of the lemma follows. \square

Remark 5.2 Note that Lemma 5.1 shows that accurate numerical integration of boundary layer functions is possible without constructing special, “exponentially” weighted quadrature rules. The present approach works even without explicit knowledge of the boundary layer function.

As we pointed out in Remark 3.4, the value of κ_0 is in principle available from the proof. For the special weight function $w = e^{-(1-x)/\varepsilon}$, the analysis of [15] shows that $\kappa_0 \geq 4/e$. Furthermore, note that the use of geometrically refined meshes outlined in Remark 4.4 eliminates the need for bounds on κ_0 . We chose π_{2q-1} as the product of (piecewise) polynomials of degree $q - 1$ and q . However, other “splittings” are possible and thus the condition $q \geq p + 1$ for polynomial right hand sides f may be relaxed to a condition of the form $q \geq \tau p$ with $\tau > 1/2$.

In Figs. 1–4, we present the results of calculations with very large values of q corresponding to practically exact evaluation of the stiffness matrix and load vector. In Figs. 1, 2 we show the L^2 convergence versus the polynomial degree (note: $\dim S_0^{\bar{p},1}(\mathcal{T}_{1,\varepsilon}) = 2p - 1$ typically). As predicted by Theorem 4.3 we have robust exponential convergence: for small values of ε the error curves are practically on top of each other. For our two-element meshes, Theorem 4.3 also gives exponential convergence of the point value at the one internal node (at $1 - \varepsilon p$). Inspection of the error at that point shows superconvergence: nodal exactness (up to machine precision) is obtained for all values of p and ε . To demonstrate the robustness of the method in the L^∞ norm is the objective of the experiments reported in Figs 3, 4. Here, the discrete L^∞ error is defined by the maximum error in sampling points

which are chosen as follows. Ω is subdivided into three sampling windows $(-1, 1 - 11p\varepsilon)$, $(1 - 11p\varepsilon, 1 - p\varepsilon)$, $(1 - p\varepsilon, 1)$ and in each window 10^4 sampling points are uniformly distributed. The graphs indicate that also in the maximum norm, the finite element solution features indeed the robust exponential convergence predicted in Theorem 3.3.

Next, we address the effect of approximate test functions on the performance of the method. By reducing the order of integration q we introduce into the stiffness matrix and the load vectors errors which correspond to the effect of approximate test functions (which, being piecewise polynomials, would be integrated exactly). Figs. 5–8 show the effect of inexact integration. The numerical integrations of exponentials times Legendre polynomials were now performed with composite Gaussian rules of order $q = 1$, $q = p/4 + 1$, $q = p/2 + 1$, $q = 3/4p + 1$, and $q = p + 1$. In Figs. 5, 6, we show the error for “exactly integrated” load vectors, but low integration order in the flux jumps of the stiffness matrix for $\varepsilon = 10^{-5}$, $\varepsilon = 10^{-10}$. We observe that even with severe underintegration, $q = 1$, practically no instability occurs, but a consistency error of size $O(\sqrt{\varepsilon})$ is introduced. A similar phenomenon is observed for the quadrature errors of the right hand side, shown in Figs. 7, 8, where now the stiffness matrix has been integrated “exactly”. Here, underintegration (e.g., $q = 1$) leads to a saturation at an error level of roughly $O(\varepsilon)$. Despite these consistency errors the stability is not compromised, even by severe underintegration with $q = 1$, which might explain the success of h -version schemes based on often very crude, analytical approximations of the upwinded test functions. It appears, however, that, in order to avoid the $O(\sqrt{\varepsilon})$ and $O(\varepsilon)$ saturation errors observed in Figs. 5–8, one has indeed to increase the quadrature order (resp. the polynomial degree q of the approximate test functions) in accordance with Corollary 4.11, i.e., proportional to $\max(p, |\ln \varepsilon| + \ln N)$.

Let us finally comment on the sparsity pattern and the solution of the resulting linear system. If the basis of the L -splines is chosen as in (3.8), (3.9), and if the basis of the trial space is the usual “nodal” and “internal” shape functions, then the resulting stiffness matrix is a banded matrix with bandwidth $O(p)$ (cf. Fig. 9 for the case of a four-element mesh and $p = 10$, i.e., 39 unknowns).

In summary, our numerical experiments show that our error estimates are sharp and that they describe accurately the performance of the Petrov-Galerkin hp -FEM: the impact of the quadrature order on the stability and consistency follows closely the predictions made in Corollary 4.11 and the method performs uniformly well for the viscosity parameter ε ranging from $\varepsilon = 1$ to the order of machine precision, $\varepsilon = 10^{-16}$. In the latter case, within the machine precision the hyperbolic limiting problem is calculated. Thus, the method presented here also opens new avenues to generate p and hp version FEM for hyperbolic problems via a numerical vanishing viscosity approach.

A Appendix: Analysis of Petrov-Galerkin FEM

Here we present some abstract results that are used repeatedly in our analysis. We merely cite those that are classical (see [1]), and derive some extensions required by us.

Throughout this appendix, X and Y denote reflexive Banach spaces equipped with norms $\|\cdot\|_X$ and $\|\cdot\|_Y$, respectively, and $B : X \times Y \rightarrow \mathbb{R}$ denotes a bilinear form which is

continuous, i.e.

$$|B(u, v)| \leq C_1 \|u\|_X \|v\|_Y. \quad (\text{A.1})$$

By Y' we denote the dual space to Y . It is also a Banach space equipped with the norm

$$\|F\|_{Y'} = \sup_{0 \neq v \in Y} \frac{|F(v)|}{\|v\|_Y}. \quad (\text{A.2})$$

Proposition A.1 *Assume that $B(\cdot, \cdot)$ satisfies*

$$\inf_{0 \neq u \in X} \sup_{0 \neq v \in Y} \frac{B(u, v)}{\|u\|_X \|v\|_Y} \geq \gamma > 0 \quad (\text{A.3})$$

and

$$\forall 0 \neq v \in Y : \sup_{u \in X} B(u, v) > 0. \quad (\text{A.4})$$

Then, for every $F \in Y'$, the abstract saddle point problem:

$$u \in X : \quad B(u, v) = F(v) \quad \forall v \in Y \quad (\text{A.5})$$

admits a unique solution u satisfying the a-priori estimate

$$\|u\|_X \leq \frac{C_1}{\gamma} \|F\|_{Y'}. \quad (\text{A.6})$$

For a proof, we refer to [1].

An equivalent form of the stability condition (A.3), (A.4) is

Proposition A.2 *If $B(\cdot, \cdot)$ satisfies*

$$\inf_{0 \neq v \in Y} \sup_{0 \neq u \in X} \frac{B(u, v)}{\|u\|_X \|v\|_Y} \geq \tilde{\gamma} > 0 \quad (\text{A.7})$$

and

$$\forall 0 \neq u \in X : \sup_{v \in Y} B(u, v) > 0 \quad (\text{A.8})$$

then also (A.3) and (A.4) hold with $\gamma = \tilde{\gamma}/C_1$.

Proof: Let $G \in X'$ satisfy $\|G\|_{X'} = 1$ and consider the auxiliary problem:

$$\tilde{v}_G \in Y : \quad B(w, \tilde{v}_G) = G(w) \quad \forall w \in X. \quad (\text{A.9})$$

Clearly, the bilinear form $C(\cdot, \cdot)$ defined via $C(v, u) := B(u, v)$ satisfies by our assumptions (A.7), (A.8) all requirements for Proposition A.1 with X and Y interchanged, however. Hence \tilde{v}_G exists, is unique and satisfies

$$\|\tilde{v}_G\|_Y \leq \frac{C_1}{\tilde{\gamma}} \|G\|_{X'} = \frac{C_1}{\tilde{\gamma}}. \quad (\text{A.10})$$

We prove that $B(\cdot, \cdot)$ satisfies (A.3): given $0 \neq u \in X$, select $G(\cdot) \in X'$ such that $G(u) = 1$ and define $v_u := \|u\|_X \tilde{v}_G$. Then, by (A.9),

$$B(u, v_u) = \|u\|_X B(u, \tilde{v}_G) = \|u\|_X^2,$$

and, by (A.10),

$$\|v_u\|_Y \leq \frac{C_1}{\gamma} \|u\|_X.$$

This implies (A.3) with $\gamma = \tilde{\gamma}/C_1$. (A.4) follows directly from (A.7). \square

Let now $X_M \subset X$, $Y_M \subset Y$ be closed subspaces. We consider the abstract FE-discretization of (A.5)

$$u_M \in X_M : \quad B(u_M, v) = F(v) \quad \forall v \in Y_M. \quad (\text{A.11})$$

The following result, due to Babuška [1], addresses the convergence of (A.11) in terms of the approximability of u from X_M and in terms of the stability implied by the test function space Y_M .

Proposition A.3 *Assume*

$$\inf_{0 \neq u \in X_M} \sup_{0 \neq v \in Y_M} \frac{B(u, v)}{\|u\|_X \|v\|_Y} \geq \gamma_M > 0 \quad (\text{A.12})$$

and

$$\forall 0 \neq v \in Y_M : \sup_{u \in X_M} B(u, v) > 0. \quad (\text{A.13})$$

Then, for every $F \in Y'$, (A.11) admits a unique solution u_M which satisfies the error estimate

$$\|u - u_M\|_X \leq \left(1 + \frac{C_1}{\gamma_M}\right) \inf_{w \in X_M} \|u - w\|_X. \quad (\text{A.14})$$

Frequently in this paper, one does not have the inf-sup conditions (A.12), (A.13), but rather the adjoint set of conditions

$$\inf_{0 \neq v \in Y_M} \sup_{0 \neq u \in X_M} \frac{B(u, v)}{\|u\|_X \|v\|_Y} \geq \tilde{\gamma}_M > 0, \quad (\text{A.15})$$

and

$$\forall 0 \neq u \in X_M : \sup_{v \in Y_M} B(u, v) > 0. \quad (\text{A.16})$$

An application of Proposition A.2 to the finite dimensional case gives

Proposition A.4 *Assume (A.1), (A.15), (A.16). Then the inf-sup conditions (A.12), (A.13) hold with*

$$\gamma_M = \tilde{\gamma}_M / C_1.$$

B Appendix: Regularity

The goal of this subsection is to prove Theorem 1.1, i.e., obtain bounds on the u_ε and its derivatives which depend *only* on the constants $C_a, C_b, C_f, \gamma_a, \gamma_b, \gamma_f$, and γ_1, γ_2 of Section 1.2

We introduce two more expressions λ^-, λ^+ which can be controlled in terms of γ_1, γ_2 :

$$\begin{aligned}\lambda^- &:= \max \left\{ \frac{a - \sqrt{a^2 + 4b\varepsilon}}{2\varepsilon}, 0 \right\} \leq \gamma_2, \\ \lambda^+ &:= \min \left\{ \frac{a + \sqrt{a^2 + 4b\varepsilon}}{2\varepsilon}, \frac{a}{\varepsilon} \right\} \geq \frac{a}{2\varepsilon}.\end{aligned}$$

Let us first get bounds on the solution u_ε by the maximum principle.

Lemma B.1 *There is $C > 0$ depending only on the constants appearing in (1.5), (1.6), (1.3), and C_f such that the solution u_ε of (1.1) satisfies*

$$\|u_\varepsilon\|_{L^\infty} \leq C, \quad \|u'_\varepsilon\|_{L^\infty} \leq C\varepsilon^{-1}.$$

Proof: The lemma is proved by the maximum principle (cf. [13], Chap. 1, Sec. 5, Thm. 11; note that the function $w := e^{-\lambda^+(1-x)}$ satisfies the assumptions of Thm. 11). Consider the functions

$$\psi_\pm := |\alpha^-|e^{\lambda^-(1+x)} + |\alpha^+|e^{-\lambda^+(1-x)} + \frac{x+1}{\gamma_1}\|f\|_{L^\infty}e^{\lambda^-(1+x)} \pm u_\varepsilon.$$

Then $\psi_\pm(\pm 1) \geq 0$, $L_\varepsilon\psi_\pm \geq 0$ and thus by the maximum principle

$$|u_\varepsilon(x)| \leq |\alpha^-|e^{\lambda^-(1+x)} + |\alpha^+|e^{-\lambda^+(1-x)} + \frac{x+1}{\gamma_1}\|f\|_{L^\infty}e^{\lambda^-(1+x)}.$$

The bound on $\|u_\varepsilon\|_{L^\infty}$ follows. Let us introduce the shorthand

$$A(x) := \frac{1}{\varepsilon} \int_x^1 a(t) dt.$$

For the derivative estimate, we estimate $u'_\varepsilon(1)$ first. Multiplying the differential equation by $e^{A(x)}$, then integrating from x to 1 and then multiplying again by $e^{-A(x)}$ gives

$$u'_\varepsilon(x) = e^{-A(x)}u'_\varepsilon(1) - \frac{1}{\varepsilon} \int_x^1 b(t)e^{A(t)-A(x)}u_\varepsilon(t) dt + \frac{1}{\varepsilon} \int_x^1 f(t)e^{A(t)-A(x)} dt \quad (\text{B.1})$$

Integrating this equation from -1 to 1 yields

$$\alpha^+ - \alpha^- = u'_\varepsilon(1) \int_{-1}^1 e^{-A(x)} dx - \frac{1}{\varepsilon} \int_{-1}^1 e^{-A(x)} \int_x^1 b(t)e^{A(t)}u_\varepsilon(t) dt dx + \frac{1}{\varepsilon} \int_{-1}^1 e^{-A(x)} \int_x^1 f(t)e^{A(t)} dt dx$$

Some simple algebra shows that we have

$$\int_{-1}^1 e^{-A(x)} dx \leq \frac{\varepsilon}{a} \quad (\text{B.2})$$

$$\int_{-1}^1 e^{-A(x)} dx \geq \frac{\varepsilon}{\|a\|_{L^\infty}}(1 - e^{2a/\varepsilon}) \quad (\text{B.3})$$

$$\frac{1}{\varepsilon} \int_{-1}^1 \int_x^1 e^{A(t)-A(x)} dt dx \leq \frac{2}{a} \quad (\text{B.4})$$

$$\frac{1}{\varepsilon} \int_x^1 e^{A(t)-A(x)} e^{-\lambda^+(1-t)} dt \leq \frac{2}{\varepsilon} e^{-\lambda^+(1-x)} \quad (\text{B.5})$$

Therefore,

$$|u'_\varepsilon(1)| \leq \left[|\alpha^+ - \alpha^-| + (\|b\|_{L^\infty} \|u_\varepsilon\|_{L^\infty} + \|f\|_{L^\infty}) \frac{2}{\underline{a}} \right] \frac{\|a\|_{L^\infty}}{\varepsilon} (1 - e^{-2\underline{a}/\varepsilon})^{-1}.$$

Thus, inserting this estimate in (B.1) yields

$$|u'_\varepsilon(x)| \leq |u'_\varepsilon(1)| + \frac{2}{\varepsilon} \|b\|_{L^\infty} \|u_\varepsilon\|_{L^\infty} + \frac{2}{\varepsilon} \|f\|_{L^\infty}$$

and thus the claim of the lemma. \square

Lemma B.2 *Let u_ε^+ be the outflow boundary layer defined in (1.9). Then there is $C > 0$ depending only on the constants of (1.3), ((1.5))–((1.7)) such that*

$$|u_\varepsilon^+(x)| \leq e^{-\underline{a}(1-x)/(2\varepsilon)}, \quad |u_\varepsilon^{+'}(x)| \leq C\varepsilon^{-1} e^{-\underline{a}(1-x)/(2\varepsilon)}.$$

Proof: In fact, a stronger statement holds true:

$$|u_\varepsilon^+(x)| \leq e^{-\lambda^+(1-x)}, \quad |u_\varepsilon^{+'}(x)| \leq C\varepsilon^{-1} e^{-\lambda^+(1-x)}.$$

The pointwise estimate follows immediately from the comparison functions used in the proof of Lemma B.1. The derivative estimate follows similarly to the one in Lemma B.1. We obtain the same bound on $u'_\varepsilon(1)$ and then insert this bound in (B.1) making use of (B.5). The observation $\lambda^+ \geq \underline{a}/(2\varepsilon)$ concludes the proof of the lemma. \square

Lemma B.3 *Let u_ε^+ be the outflow boundary layer defined in (1.9). Then there are constants $C_1, C_2 > 0$ depending only on the constants of (1.3), ((1.5)), ((1.6)) such that*

$$C_1\varepsilon^{-1} \leq u_\varepsilon^+(1) \leq C_2\varepsilon^{-1}.$$

Proof: Lemma B.2 gives the upper bound. For the lower bound, we analyze the proof of Lemma B.2 more carefully. Let the function A be defined as in the proof of Lemma B.2. The equation following (B.1) reads

$$1 = u_\varepsilon^{+'}(1) \int_{-1}^1 e^{-A(x)} dx - \frac{1}{\varepsilon} \int_{-1}^1 e^{-A(x)} \int_x^1 b(t) e^{A(t)} u_\varepsilon^+(t) dt dx \quad (\text{B.6})$$

By the maximum principle and Lemma B.2 we have

$$0 \leq u_\varepsilon^+(x) \leq e^{-\lambda^+(1-x)}, \quad x \in \overline{\Omega} \quad (\text{B.7})$$

Thus, if $b \geq \underline{b} \geq 0$ on $\overline{\Omega}$, then the claim of the Lemma follows with (B.3). Let therefore $\underline{b} < 0$. We obtain again with $u_\varepsilon^+ \geq 0$

$$u_\varepsilon^{+'}(1) \int_{-1}^1 e^{-A(x)} dx \geq 1 - \frac{1}{\varepsilon} \int_{-1}^1 e^{-A(x)} \int_x^1 |\underline{b}| e^{A(t)} u_\varepsilon^+(t) dt dx \quad (\text{B.8})$$

Using $\underline{a} \leq a(x)$, $-A'(x) = a(x)/\varepsilon$, and (B.7) we obtain

$$\begin{aligned} \frac{1}{\varepsilon} \int_{-1}^1 e^{-A(x)} \int_x^1 |\underline{b}| e^{A(t)} u_\varepsilon^+(t) dt dx &\leq \int_{-1}^1 \frac{a(x)}{\varepsilon \underline{a}} e^{-A(x)} \int_x^1 |\underline{b}| e^{A(t)} u_\varepsilon^+(t) dt dx \\ &\leq \frac{1}{\underline{a}} \left[e^{-A(x)} \int_x^1 |\underline{b}| e^{A(t)} u_\varepsilon^+(t) dt \right]_{-1}^1 + \frac{1}{\underline{a}} \int_{-1}^1 |\underline{b}| u_\varepsilon^+(t) dt \\ &\leq \frac{1}{\underline{a}} \int_{-1}^1 |\underline{b}| e^{-\lambda^+(1-t)} dt \leq \frac{|\underline{b}|}{\underline{a} \lambda^+} \leq \frac{2|\underline{b}| \varepsilon}{\underline{a}^2} \end{aligned}$$

Inserting this estimate in (B.8) and noting that $|\underline{b}| = -\underline{b}$ leads to

$$u_\varepsilon^{+'}(1) \int_{-1}^1 e^{-A(x)} dx \geq 1 - \frac{2|\underline{b}| \varepsilon}{\underline{a}^2} = (\underline{a}^2 + 2\underline{b}\varepsilon) / \underline{a}^2 \geq \gamma_1^2 \underline{a}^2$$

and thus the claim of the lemma by estimate (B.3). \square

Proof of (1.12), (1.13): Let us first prove (1.12). Choose $K > \max(1, \gamma_f, \gamma_a, \gamma_b)$ such that

$$\left[\frac{C_f}{K^2} + \frac{C_a}{K} \frac{1}{1 - \gamma_a/K} + \frac{C_b}{K^2} \frac{1}{1 - \gamma_b/K} \right] \leq 1.$$

By Lemma B.1, we may now choose the constant $C \geq 1$ such that (1.12) holds true for $n = 0, 1$. Let us now proceed by induction on n . We assume that the induction hypothesis (1.12) holds for $0 \leq \nu \leq n+1$ and show that it holds for $n+2$. Differentiating the differential equation (1.1) n times (note that we know already that u_ε is analytic) we get

$$-\varepsilon u_\varepsilon^{(n+2)} = f^{(n)} - (au'_\varepsilon)^{(n)} - (bu_\varepsilon)^{(n)} = f^{(n)} - \sum_{\nu=0}^n \binom{n}{\nu} \left(a^{(\nu)} u_\varepsilon^{(n+1-\nu)} + b^{(\nu)} u_\varepsilon^{(n-\nu)} \right).$$

Using the induction hypothesis, we get

$$\begin{aligned} \varepsilon \|u_\varepsilon^{(n+2)}\|_{L^\infty(\Omega)} &\leq \|f^{(n)}\|_{L^\infty(\Omega)} + C \sum_{\nu=0}^n \binom{n}{\nu} \left[C_a \gamma_a^\nu \nu! \max(n+1-\nu, \varepsilon^{-1})^{n+1-\nu} \right. \\ &\quad \left. + C_b \gamma_b^\nu \nu! C K^{n-\nu} \max(n-\nu, d^{-1})^{n-\nu} \right]. \end{aligned}$$

Exploiting the estimates

$$\begin{aligned} \binom{n}{\nu} \nu! \max(n+1-\nu, \varepsilon^{-1})^{n+1-\nu} &\leq n^\nu \max(n+1, \varepsilon^{-1})^{n+1-\nu} \leq \max(n+1, \varepsilon^{-1})^{n+1} \\ \binom{n}{\nu} \nu! \max(n-\nu, \varepsilon^{-1})^{n-\nu} &\leq n^\nu \max(n, \varepsilon^{-1})^{n-\nu} \leq \max(n+1, \varepsilon^{-1})^{n+1} \\ \|f^{(n)}\|_{L^\infty(\Omega)} &\leq C_f \gamma_f^n n! \leq C_f \max(n+1, \varepsilon^{-1})^{n+1} \end{aligned}$$

we obtain

$$\begin{aligned} \varepsilon \|u_\varepsilon^{(n+2)}\|_{L^\infty(\Omega)} &\leq \|f^{(n)}\|_{L^\infty(\Omega)} + C K^{n+2} \max(n+1, \varepsilon^{-1})^{n+1} \sum_{\nu=0}^n \frac{C_a}{K} \left(\frac{\gamma_a}{K}\right)^\nu + \frac{C_b}{K^2} \left(\frac{\gamma_b}{K}\right)^\nu \\ &\leq C K^{n+2} \max(n+1, \varepsilon^{-1})^{n+1} \left[\frac{C_f}{K^2} + \frac{C_a}{K} \frac{1}{1 - \gamma_a/K} + \frac{C_b}{K^2} \frac{1}{1 - \gamma_b/K} \right]. \end{aligned}$$

By the choice of K the expression in the brackets is bounded by 1 which concludes the induction argument after dividing both sides by ε .

The proof of (1.13) proceeds in the same fashion; the only difference is that Lemma B.2 instead of Lemma B.1 is used to start the induction argument.

□

Now we turn to the proof of (1.14)—(1.16). Recall the definition of the terms u_j of the asymptotic part w_M in the decomposition (1.9). In order to control these terms, we need the following lemma.

Lemma B.4 *Let G be an open, complex neighborhood of $I = [-1, 1]$. Assume that the functions $\Lambda, a, u_0 : G \rightarrow \mathbf{C}$ are holomorphic and bounded on G . Assume additionally that $|a| \geq \underline{a} > 0$ on G . Then there are constants $C, K_1, K_2 > 0$ depending only on $\underline{a}, \|a'\|_{L^\infty(G)}, \|\Lambda\|_{L^\infty(G)}, \|\Lambda'\|_{L^\infty(G)}$, and G such that the functions u_j defined recursively as in (1.8) satisfy*

$$\|u_j^{(n)}\|_{L^\infty(I)} \leq CK_1^j K_2^n j! n! \|u_0\|_{L^\infty(G)} \quad \forall j, n \in \mathbf{N}_0.$$

Proof: Again, we will prove a stronger statement. Without loss of generality we may assume that G is star shaped with respect to $z = -1$. For $\delta > 0$ (sufficiently small) denote $G_\delta := \{z \in G \mid \text{dist}(z, \partial G) > \delta\}$. Then we claim that there are $C, K > 0$ such that

$$\|u_j\|_{L^\infty(G_\delta)} \leq CK^j \delta^{-j} j! \|u_0\|_{L^\infty(G)} \quad \forall j \in \mathbf{N}_0.$$

The proof of the lemma follows from this estimate by Cauchy's integral theorem for derivatives.

It remains therefore to establish the claim. We proceed by induction on j . It is true for $j = 0$ and for any $C \geq 1$. We write

$$\begin{aligned} u_{j+1}(z) &= e^{-\Lambda(z)} \int_{-1}^z e^{\Lambda(t)} \frac{1}{a(t)} u_j''(t) dt \\ &= e^{-\Lambda(z)} \left[e^{\Lambda(t)} \frac{1}{a(t)} u_j'(t) \right]_{-1}^z - e^{-\Lambda(z)} \int_{-1}^z e^{\Lambda(t)} \frac{\Lambda'(t)a(t) + a'(t)}{a(t)^2} u_j'(t) dt. \end{aligned}$$

Hence there is $C_1 > 0$ such that

$$\|u_{j+1}\|_{L^\infty(G_\delta)} \leq C_1 \|u_j'\|_{L^\infty(G_\delta)}.$$

By Cauchy's integral theorem, we have for $0 < \kappa < 1$ using the induction hypothesis:

$$\begin{aligned} \|u_{j+1}\|_{L^\infty(G_\delta)} &\leq C_1 \frac{2\pi\kappa\delta}{(\kappa\delta)^2} \|u_j\|_{L^\infty(G_{(1-\kappa)\delta})} \\ &\leq C_1 C j! K^j (1-\kappa)^{-j} \delta^{-j} \frac{1}{\kappa\delta} \|u_0\|_{L^\infty(G)} \\ &\leq C(j+1)! K^{j+1} \delta^{-(j+1)} \|u_0\|_{L^\infty(G)} \frac{C_1}{K(j+1)(1-\kappa)^j \kappa} \end{aligned}$$

Choosing $\kappa = 1/(j + 1)$, we observe that there is $c_1 > 0$ such that $c_1 \leq (j + 1)\kappa(1 - \kappa)^j$ for all $j \geq 1$. Hence, choosing $K > 0$ such that $C_1/(Kc_1) \leq 1$ finishes the induction argument. \square

This lemma puts us in position to conclude the proof of Theorem 1.1.

Proof of (1.14)–(1.16): Let us begin with (1.14). We see immediately that the assumptions of Lemma B.4 are satisfied: The size of the complex neighborhood G and the constants C , K_1 , K_2 can be controlled by the constants of (1.3), (1.5)–(1.7). Also $\|u_0\|_{L^\infty(G)}$ can be controlled in terms of these constants. Hence, the terms u_j satisfy

$$\|u_j^{(n)}\|_{L^\infty(\Omega)} \leq Cj!n!K_1^jK_2^n \quad \forall j \in \mathbb{N}_0, n \in \mathbb{N}_0.$$

Thus

$$\|w_M^{(n)}\|_{L^\infty(\Omega)} \leq CK_2^n n! \sum_{j=0}^M \varepsilon^j K_1^j j! \leq CK_2^n n! \sum_{j=0}^M (\varepsilon K_1 M)^j.$$

This last sum can be bounded by a constant under the condition $\varepsilon KM \leq 1$ if we choose $K > K_1$.

An immediate consequence of (1.14) is (1.16): $C_M = \alpha^+ - w_M(1)$ and $w_M(1)$ can be controlled by (1.14) under the assumption $\varepsilon MK \leq 1$. Finally, the remainder r_M satisfies (1.11). An application of Lemma B.1 (note that we only need to control the L^∞ norm of the right hand side for Lemma B.1 to hold) together with (1.14) gives the claim of (1.16) for $n = 0, 1$ and the differential equation satisfied by r_M gives the claim for $n = 2$. \square

References

- [1] I. Babuška and A.K.Aziz: *Survey lectures on the mathematical foundations of the finite element method*, in: The mathematical foundations of the finite element method, A.K.Aziz and I. Babuška (Eds.), Academic Press 1972.
- [2] J. Christie, D.F. Griffiths, A.R. Mitchell, and O.C. Zienkiewicz, *Finite Element Methods for Second Order Differential Equations with Significant First Derivatives*, Int. J. Num. Meths. in Eng. **10** (1978), 1764–1771
- [3] P.J. Davis, *Interpolation and Approximation*, Dover Publ. 1974.
- [4] E.C. Gartland, *Uniform high-order difference schemes for a singularly perturbed two-point boundary value problem*, Math. Comp. **48** 551-564 (1987).
- [5] P.P.N. De Groen, *A finite element method with large mesh width for stiff two-point boundary value problems* Preprint, Dept. Mathematics, Eindhoven Univ., Eindhoven, The Netherlands, 1978.

- [6] P.P.N. De Groen and P.W. Hemker, *Error bounds for exponentially fitted Galerkin methods applied to stiff two-point boundary value problems*, in Numerical Analysis of Singular Perturbation Problems, P.W. Hemker and J.J.H. Miller, eds., Academic Press, New York, 1979
- [7] P.W. Hemker, *A Numerical Study of Stiff Two-Point Boundary Problems*, PhD thesis, Mathematisch Centrum, Amsterdam, 1977
- [8] A.I. Markushevich, *Theory of Functions of a Complex Variable*, Chelsea Publishing Company, 1977.
- [9] J.M. Melenk, *On robust exponential convergence of finite element methods for problems with boundary layers*, Research Report 96-06, Sem. for Appl. Math. ETH Zürich, May 1996 (submitted to IMA J. Num. Anal.).
- [10] J.M. Melenk and C. Schwab, in preparation
- [11] Miller, J.J.H, O’Riordan, E., and Shishkin, G.I., *Fitted Numerical Methods for Singular Perturbation Problems*, World Scientific Publishers, Singapore 1996.
- [12] K.W. Morton, *Numerical solution of convection-diffusion problems*, Oxford Univ. Press, Oxford 1995.
- [13] M. Protter and H. Weinberger, *Maximum Principles in Differential Equations*, Springer Verlag Heidelberg-New York 1984.
- [14] H.G. Roos, M. Stynes, and L. Tobiska, *Numerical Solution of singularly perturbed boundary value problems*, Springer Verlag Heidelberg, New York 1995.
- [15] C. Schwab and M. Suri, *The p and hp versions of the finite element method for problems with boundary layers*, Math. Comp. **65** 1403–1429 (1996).
- [16] C. Schwab, M. Suri, and C.A. Xenophontos, *Boundary Layer Approximation by Spectral/ hp Methods*, Houston J. Math. Spec. Issue of ICOSAHOM ’95 Conference A.V. Illin and L.R. Scott (Eds.), 501-508 (1996) .
- [17] B. Sündermann *Lebesgue constants in Lagrangian interpolation at the Fekete points*, Ergebnisberichte der Lehrstühle Mathematik III und VIII (Angewandte Mathematik) 44, Universität Dortmund, 1980.
- [18] G.W. Szymczak, *An adaptive Finite Element Method for convection-diffusion problems*, Ph.D. Dissertation Univ. Maryland College Park 1982.
- [19] G.W. Szymczak and I. Babuška, *Adaptivity and error estimation for the finite element method applied to convection diffusion problems* SIAM J. Num. Anal. **21** 910-954 (1984).

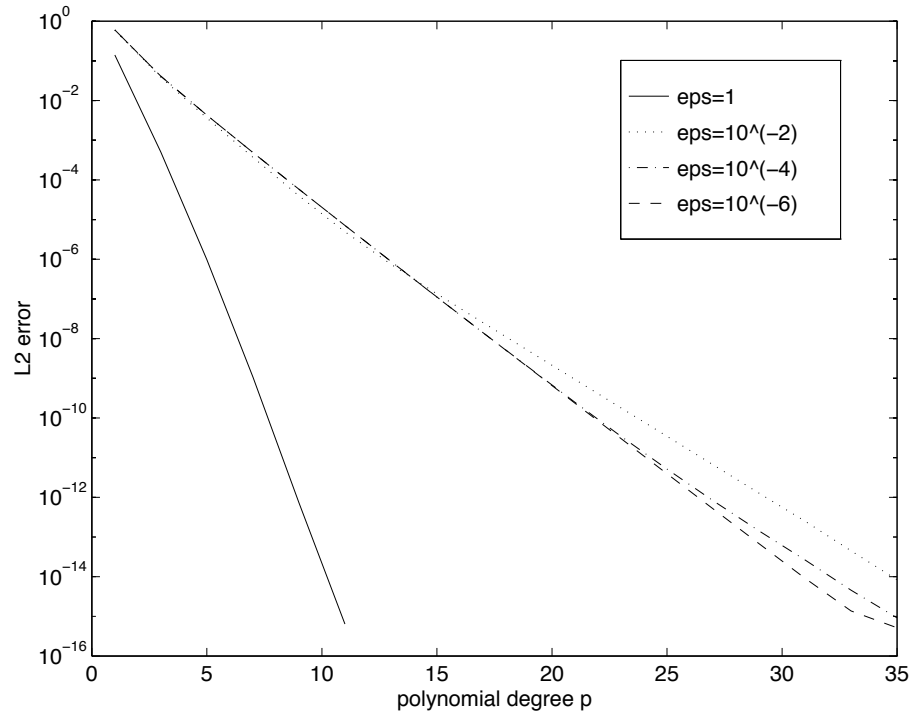


Figure 1: L^2 performance of "two-element mesh"

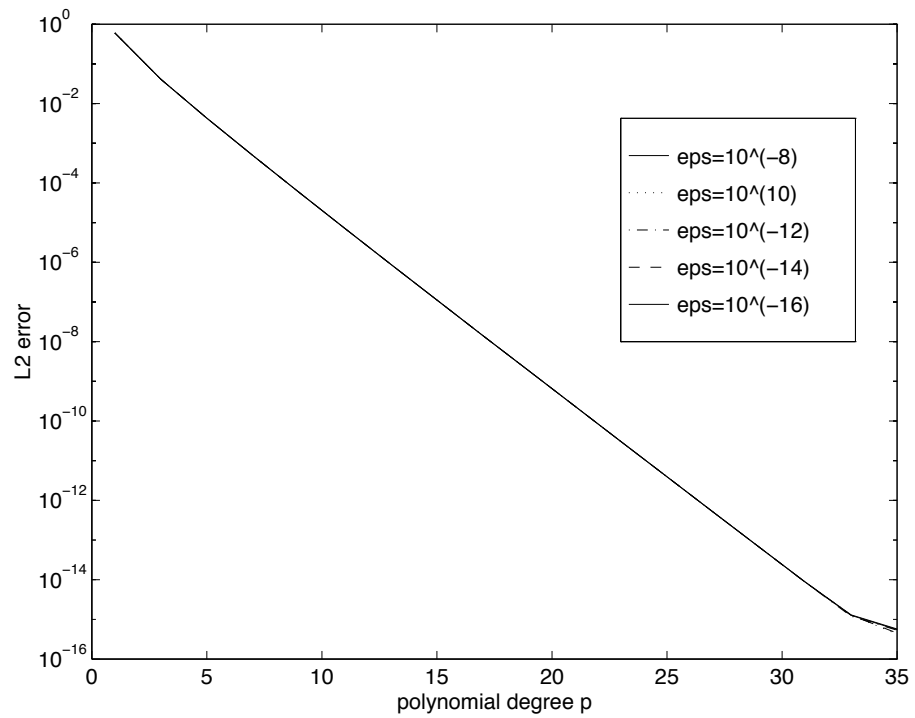


Figure 2: L^2 performance of "two-element mesh"

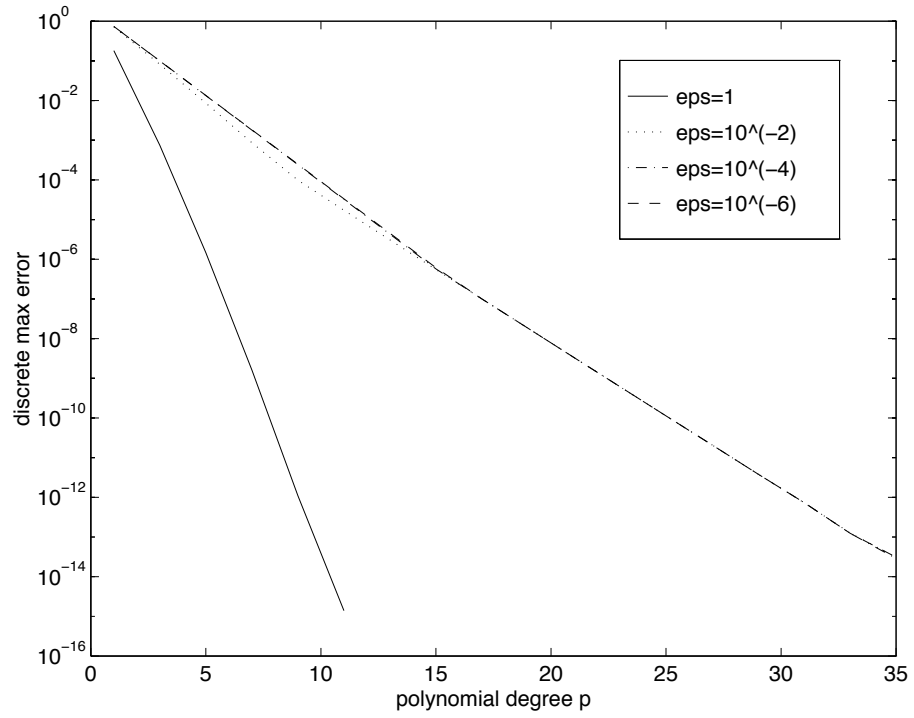


Figure 3: Performance of "two-element mesh" for discrete L^∞ norm

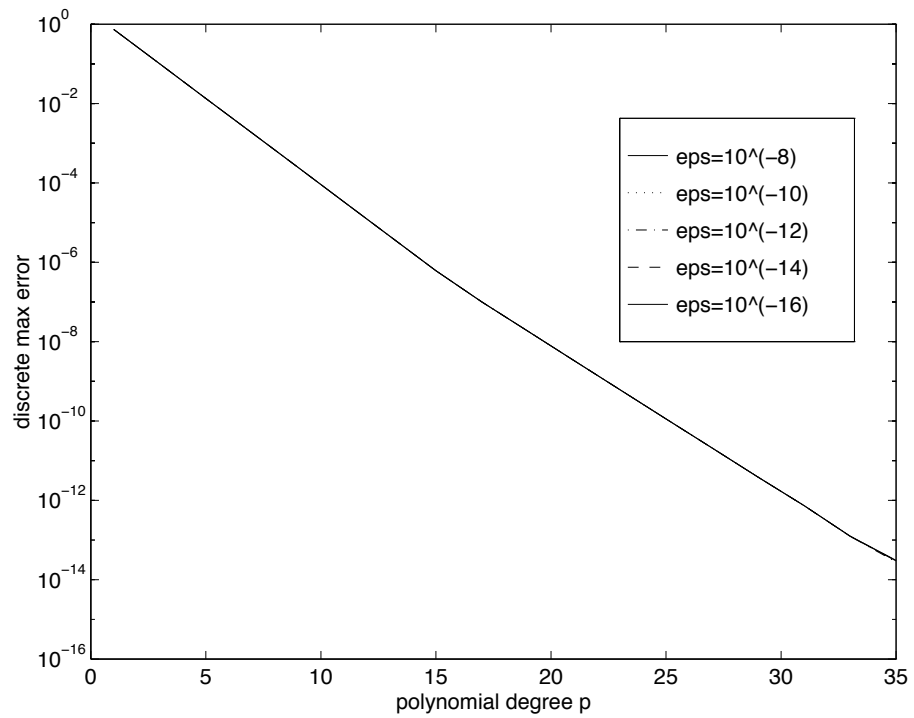


Figure 4: Performance of "two-element mesh" for discrete L^∞ norm

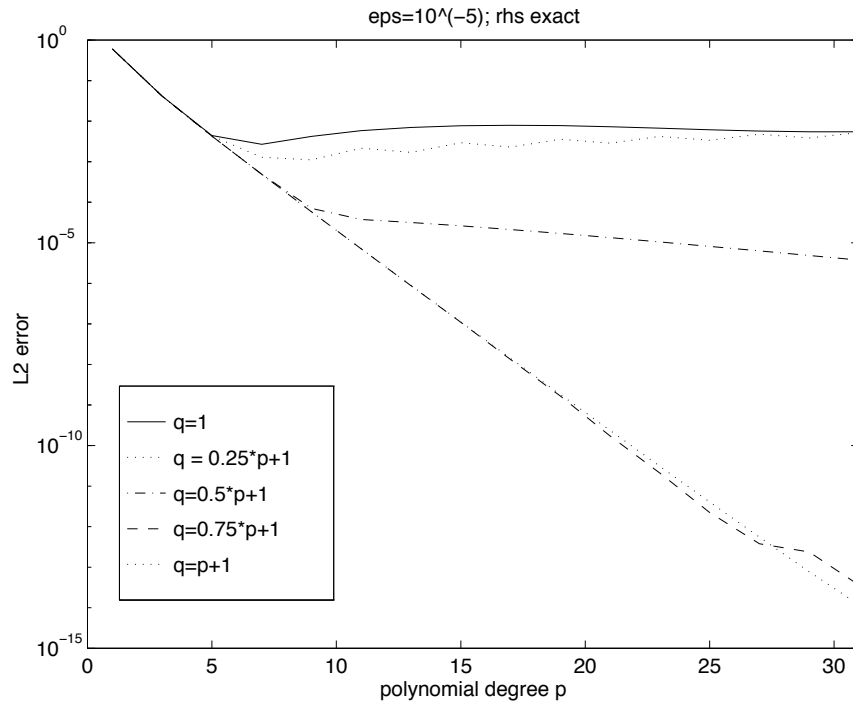


Figure 5: Effect of numerical integration of stiffness matrix; $\varepsilon = 10^{-5}$

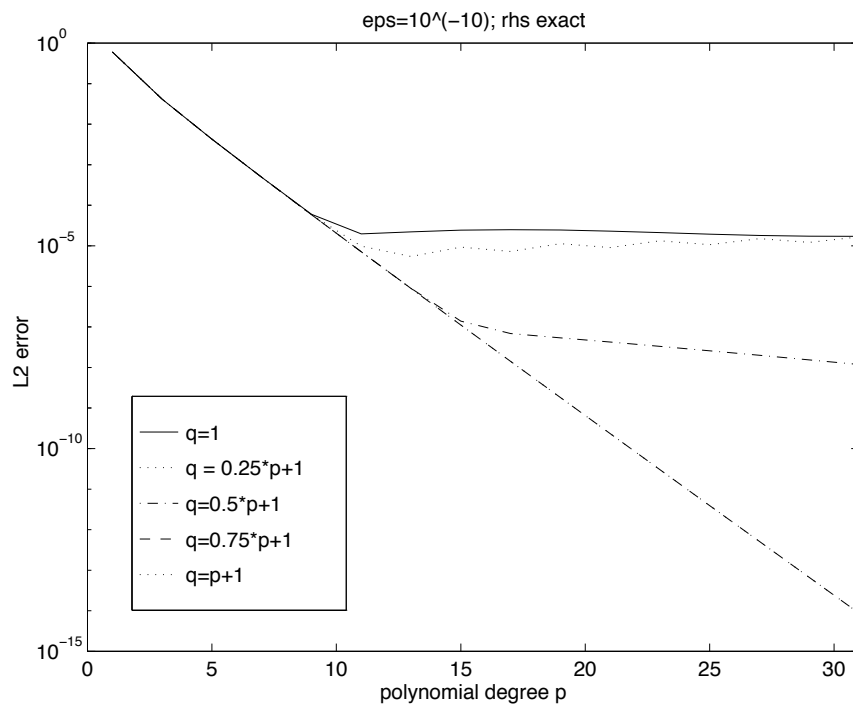


Figure 6: Effect of numerical integration of stiffness matrix; $\varepsilon = 10^{-10}$

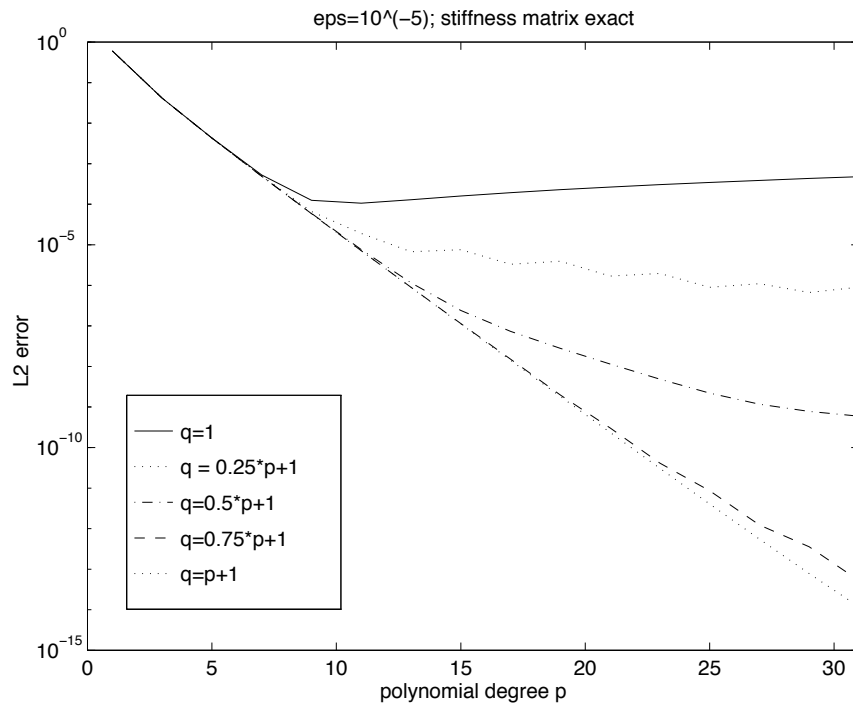


Figure 7: Effect of numerical integration of right hand side; $\varepsilon = 10^{-5}$

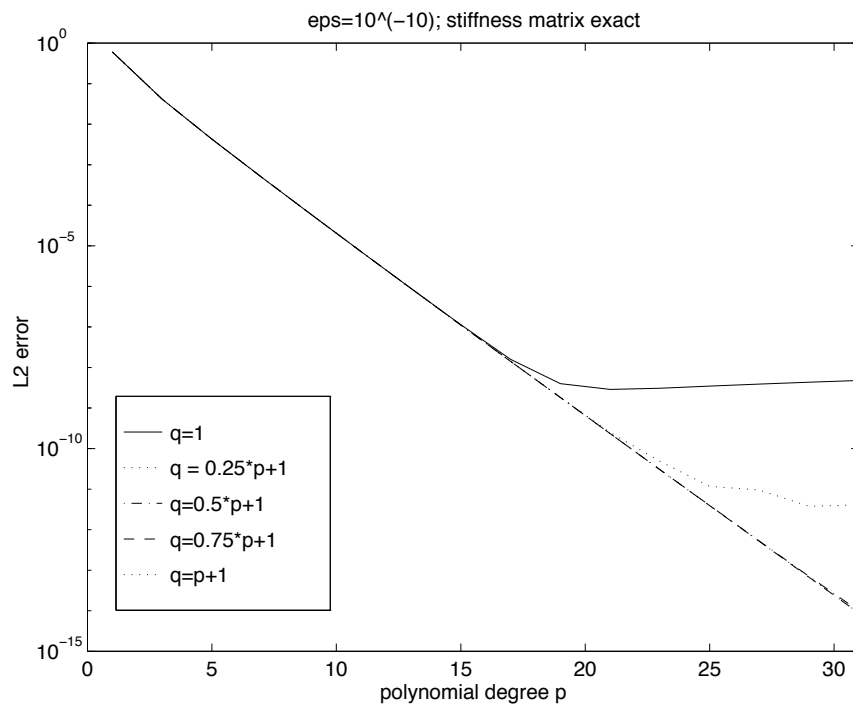


Figure 8: Effect of numerical integration of right hand side; $\varepsilon = 10^{-10}$

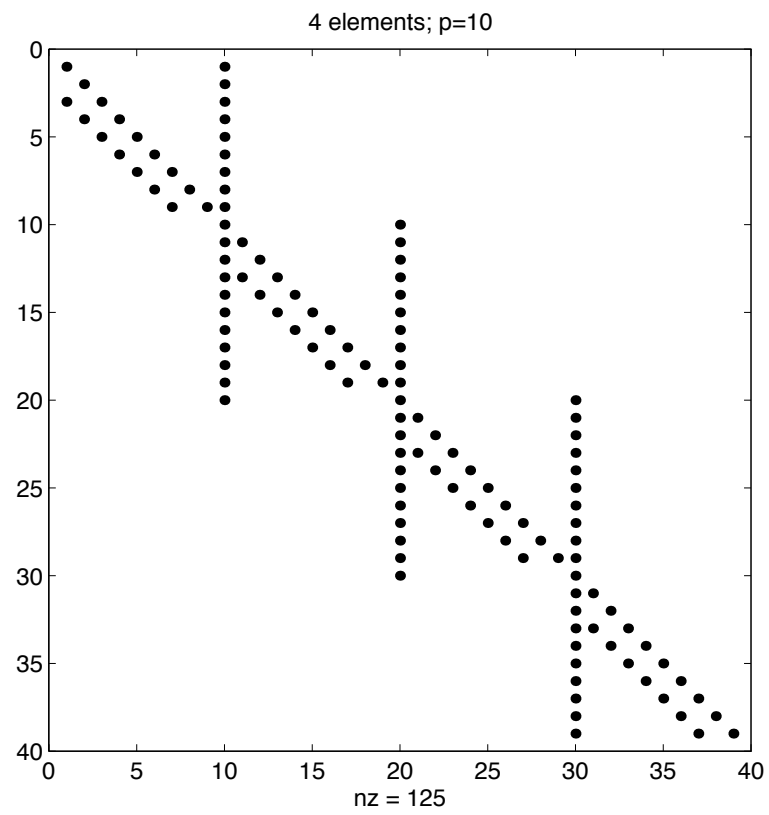


Figure 9: Sparsity structure of the stiffness matrix for four element mesh and $p = 10$

Research Reports

No.	Authors	Title
97-05	J.M. Melenk, C. Schwab	An <i>hp</i> Finite Element Method for convection-diffusion problems
97-04	J.M. Melenk, C. Schwab	<i>hp</i> FEM for Reaction-Diffusion Equations. II. Regularity Theory
97-03	J.M. Melenk, C. Schwab	<i>hp</i> FEM for Reaction-Diffusion Equations. I: Robust Exponential Convergence
97-02	D. Schötzau, C. Schwab	Mixed <i>hp</i> -FEM on anisotropic meshes
97-01	R. Sperb	Extension of two inequalities of Payne
96-22	R. Bodenmann, A.-T. Morel	Stability analysis for the method of transport
96-21	K. Gerdes	Solution of the 3D-Helmholtz equation in exterior domains of arbitrary shape using <i>HP</i> -finite infinite elements
96-20	C. Schwab, M. Suri, C. Xenophontos	The <i>hp</i> finite element method for problems in mechanics with boundary layers
96-19	C. Lage	The Application of Object Oriented Methods to Boundary Elements
96-18	R. Sperb	An alternative to Ewald sums. Part I: Identities for sums
96-17	M.D. Buhmann, Ch.A. Micchelli, A. Ron	Asymptotically Optimal Approximation and Numerical Solutions of Differential Equations
96-16	M.D. Buhmann, R. Fletcher	M.J.D. Powell's work in univariate and multivariate approximation theory and his contribution to optimization
96-15	W. Gautschi, J. Waldvogel	Contour Plots of Analytic Functions
96-14	R. Resch, F. Stenger, J. Waldvogel	Functional Equations Related to the Iteration of Functions
96-13	H. Forrer	Second Order Accurate Boundary Treatment for Cartesian Grid Methods
96-12	K. Gerdes, C. Schwab	Hierarchic models of Helmholtz problems on thin domains
96-11	K. Gerdes	The conjugated vs. the unconjugated infinite element method for the Helmholtz equation in exterior domains
96-10	J. Waldvogel	Symplectic Integrators for Hill's Lunar Problem
96-09	A.-T. Morel, M. Fey, J. Maurer	Multidimensional High Order Method of Transport for the Shallow Water Equations
96-08	A.-T. Morel	Multidimensional Scheme for the Shallow Water Equations