

Optimization of plane wave directions in plane wave discontinuous Galerkin methods for the Helmholtz equation

Akshay Agrawal and Ronald H. W. Hoppe

(Communicated by Luís Nunes Vicente)

Abstract. Recently, the use of special local test functions other than polynomials in Discontinuous Galerkin (DG) approaches has attracted a lot of attention and became known as DG-Trefftz methods. In particular, for the 2D Helmholtz equation plane waves have been used in [11] to derive an Interior Penalty (IP) type Plane Wave DG (PWDG) method and to provide an a priori error analysis of its p -version with respect to equidistributed plane wave directions. The dependence on the distribution of the plane wave directions has been studied in [1] based on a least squares method. However, the emphasis in [1] has been on the h -version of the PWDG method, i.e., decreasing the mesh width h for a fixed number p of plane wave directions. In this contribution, we are interested in the p -version, i.e., increasing p for a fixed mesh-width h . We formulate the choice of the plane wave directions as a control constrained optimal control problem with a continuously differentiable objective functional and the variational formulation of the PWDG method as a further constraint. The necessary optimality conditions are derived and numerically solved by a projected gradient method. Numerical results are given which illustrate the benefits of the approach.

Mathematics Subject Classification: 65N30, 35J20, 74J20

Keywords: Plane Wave Discontinuous Galerkin methods, optimization of plane wave directions, Helmholtz equation

1. Introduction

The use of plane waves in the finite element approximation of the Helmholtz equation goes back to the ultra weak variational formulation of the problem by Cessenat and Després [4]. The approach can be interpreted as a Discontinuous Galerkin (DG) approximation and is therefore referred to as the Plane Wave Discontinuous Galerkin (PWDG) method. Since it uses local trial spaces consisting of plane waves, it is also a particular example of a Trefftz-type finite element approximation and hence called a Trefftz-type DG method. Due to its superior

performance compared to standard finite element approximations which suffer from the so-called pollution effect, it has been studied extensively in the literature (cf., e.g., [2], [3], [5], [7], [8], [10]). In particular, the h -version and the p -version of the PWDG method have been analyzed in [9] and in [11], whereas the exponential convergence of the hp -version has been established in [12].

The PWDG method features a triangulation $\mathcal{T}_h(\Omega)$ of the computational domain $\Omega \subset \mathbb{R}^2$ with $\text{card}(\mathcal{T}_h(\Omega)) = N$ and the use of a certain number $p = 2m + 1$, $m \in \mathbb{N}$, of plane waves in each element $K \in \mathcal{T}_h(\Omega)$ which compose the local trial spaces. The plane waves are of the form $\exp(i\omega \mathbf{d}_{(j-1)p+\ell} \cdot \mathbf{x})$, where

$$\mathbf{d}_{(j-1)p+\ell} = (\cos(\theta_{(j-1)p+\ell}), \sin(\theta_{(j-1)p+\ell}))^T, \quad 1 \leq j \leq N, 1 \leq \ell \leq p,$$

$\mathbf{x} \in K$, and ω stands for the wavenumber. It is known from the convergence analysis of the PWDG method [11] that the p directions $\mathbf{d}_{(j-1)p+\ell}$, $1 \leq \ell \leq p$, $j \in \{1, \dots, N\}$, should be chosen in such a way that the minimum angle between two different directions is greater or equal $2\pi\eta/p$ for some $\eta \in (0, 1]$. The issue how to choose the directions in order to minimize a given objective functional has been considered in [1] based on a least squares method similar to that in [15]. The emphasis in [1] has been on the h -version of the PWDG approach, i.e., decreasing the mesh width h for a fixed number p of plane wave directions. Moreover, no constraints on the directions have been observed albeit it is known from [11] that the approximation properties of the PWGD method get lost, if the difference between two different directions becomes too small.

In this paper, we are interested in the p -version, i.e., increasing p for a fixed mesh width h , and we formulate this problem as a control constrained optimal control problem with a continuously differentiable objective functional and the variational formulation of the PWDG method as a further constraint, where the controls are the Np angles $\theta_{(j-1)p+\ell}$, $0 \leq \ell \leq p$, $1 \leq j \leq N$. We derive the first order necessary optimality conditions by means of the Lagrange multiplier approach and derive a projected gradient type method with Armijo line search to compute an optimal solution. Numerical results illustrate the dependence of the L^2 -norm of the global discretization error on the choice of the plane wave directions.

2. The PWDG method

For a bounded convex polygonal domain $\Omega \subset \mathbb{R}^2$ with boundary $\Gamma = \partial\Omega$ we consider the Helmholtz equation

$$-\Delta u - \omega^2 u = 0 \quad \text{in } \Omega, \quad (2.1a)$$

$$\mathbf{n} \cdot \nabla u + i\omega u = g \quad \text{on } \Gamma = \partial\Omega. \quad (2.1b)$$

where $\omega > 0$ is the wavenumber, $g \in L^2(\Gamma)$ is a given function, and \mathbf{n} denotes the exterior unit normal vector on Γ . We rewrite (2.1) as the first order system:

$$i\omega\boldsymbol{\sigma} - \nabla u = \mathbf{0} \quad \text{in } \Omega, \quad (2.2a)$$

$$-\nabla \cdot \boldsymbol{\sigma} + i\omega u = 0 \quad \text{in } \Omega, \quad (2.2b)$$

$$i\omega \mathbf{n} \cdot \boldsymbol{\sigma} + i\omega u = g \quad \text{on } \Gamma. \quad (2.2c)$$

The variational formulation of (2.2) reads: Find $(\boldsymbol{\sigma}, u) \in \mathbf{H}(\text{div}, \Omega) \times H^1(\Omega)$ such that for all $(\boldsymbol{\tau}, v) \in \mathbf{H}(\text{div}, \Omega) \times H^1(\Omega)$ it holds

$$(i\omega\boldsymbol{\sigma}, \boldsymbol{\tau})_{0,\Omega} + (u, \nabla \cdot \boldsymbol{\tau})_{0,\Omega} = \langle u, \mathbf{n} \cdot \boldsymbol{\tau} \rangle_{H^{1/2}(\Gamma), H^{-1/2}(\Gamma)}, \quad (2.3a)$$

$$(\boldsymbol{\sigma}, \nabla v)_{0,\Omega} + (u, v)_{0,\Gamma} + (i\omega u, v)_{0,\Omega} = \left(\frac{1}{i\omega} g, v \right)_{0,\Gamma}. \quad (2.3b)$$

We consider a shape regular family of geometrically conforming, quasi-uniform simplicial triangulations $\mathcal{T}_h(\Omega)$ of the computational domain Ω . For $D \subset \overline{\Omega}$ we denote by $\mathcal{E}_h(D)$ the set of edges of the triangulation in D . For $T \in \mathcal{T}_h(\Omega)$, we refer to h_T as the diameter of T and set $h := \max\{h_T \mid T \in \mathcal{T}_h(\Omega)\}$. For $E \in \mathcal{E}_h(\overline{\Omega})$, the length of E will be denoted by h_E . For functions $v \in \prod_{T \in \mathcal{T}_h(\Omega)} H^1(T)$ the trace of v on $E \in \mathcal{E}_h(\Omega)$ may exhibit a jump across E . For $E \in \mathcal{E}_h(\Omega)$ with $E = T_+ \cap T_-$, $T_\pm \in \mathcal{T}_h(\Omega)$ and $E \in \mathcal{E}_h(\Gamma)$ we define

$$\begin{aligned} \{v\}_E &:= \begin{cases} (v|_{T_+ \cap E} + v|_{T_- \cap E})/2, & E \in \mathcal{E}_h(\Omega) \\ v|_E, & E \in \mathcal{E}_h(\Gamma) \end{cases}, \\ [v]_E &:= \begin{cases} v|_{T_+ \cap E} - v|_{T_- \cap E}, & E \in \mathcal{E}_h(\Omega) \\ v|_E, & E \in \mathcal{E}_h(\Gamma) \end{cases}. \end{aligned} \quad (2.4)$$

For vector-valued functions we use an analogous notation.

We approximate (2.3a), (2.3b) by introducing the following local spaces spanned by plane waves

$$V_p(T_j) := \left\{ v(\mathbf{x}) := \sum_{\ell=1}^p \alpha_{(j-1)p+\ell} \exp(i\omega \mathbf{d}_{(j-1)p+\ell} \cdot \mathbf{x}) \right\}, \quad p \in \mathbb{N}, \quad (2.5)$$

$$\mathbf{V}_p(T_j) := V_p(T_j)^2, \quad j \in \{1, \dots, N\},$$

where $\alpha_{(j-1)p+\ell} \in \mathbb{C}$ and $\mathbf{d}_{(j-1)p+\ell}$, $1 \leq \ell \leq p$, $j \in \{1, \dots, N\}$, are p different unit directions

$$\mathbf{d}_{(j-1)p+\ell} = (\cos(\theta_{(j-1)p+\ell}), \sin(\theta_{(j-1)p+\ell}))^T, \quad 1 \leq \ell \leq p, m \in \mathbb{N}. \quad (2.6)$$

We define $\boldsymbol{\theta} = (\theta_1, \dots, \theta_{Np})^T$ such that $0 \leq \theta_{(j-1)p+\ell} < 2\pi$, $1 \leq \ell \leq p$, for $j \in \{1, \dots, N\}$. Setting

$$\hat{\theta}_{(j-1)p+\ell} = \begin{cases} \theta_{(j-1)p+\ell}, & 1 \leq \ell \leq p \\ \theta_{(j-1)p+1} + 2\pi, & \ell = p+1, \end{cases} \quad 1 \leq j \leq N, \quad (2.7)$$

we require that $\boldsymbol{\theta} \in \mathbf{K}$, where \mathbf{K} is given as follows

$$\begin{aligned} \mathbf{K} &:= \{\boldsymbol{\theta} \in [0, 2\pi)^{Np} \mid \theta_{\min} \leq \hat{\theta}_{(j-1)p+\ell+1} - \hat{\theta}_{(j-1)p+\ell} \leq \theta_{\max}, \\ &\quad 1 \leq j \leq N, 1 \leq \ell \leq p\}, \\ \theta_{\min} &:= (2\pi\eta_1)/p, \quad \theta_{\max} := (2\pi\eta_2)/p, \quad 0 < \eta_1 < 1 < \eta_2 < 3/2. \end{aligned} \quad (2.8)$$

The associated global spaces are given by

$$\begin{aligned} V_h &:= \{v_h \in L^2(\Omega) \mid v_h|_{T_j} \in V_p(T_j), 1 \leq j \leq N\}, \\ \mathbf{V}_h &:= \{\boldsymbol{\tau}_h \in L^2(\Omega)^2 \mid \boldsymbol{\tau}_h|_{T_j} \in \mathbf{V}_p(T_j), 1 \leq j \leq N\}. \end{aligned} \quad (2.9)$$

Then, the PWDG approximation of (2.1a), (2.1b) amounts to the computation of $(u_h, \boldsymbol{\sigma}_h) \in V_h \times \mathbf{V}_h$ such that for all $(v_h, \boldsymbol{\tau}_h) \in V_h \times \mathbf{V}_h$ it holds

$$\sum_{T \in \mathcal{T}_h(\Omega)} ((i\omega \boldsymbol{\sigma}_h, \boldsymbol{\tau}_h)_{0,T} + (u_h, \nabla \cdot \boldsymbol{\tau}_h)_{0,T}) - \sum_{T \in \mathcal{F}_h(\Omega)} (\hat{u}_h, \mathbf{n}_{\partial T} \cdot \boldsymbol{\tau}_h)_{0,\partial T} = 0, \quad (2.10a)$$

$$\sum_{T \in \mathcal{F}_h(\Omega)} ((\boldsymbol{\sigma}_h, \nabla v_h)_{0,T} + (i\omega u_h, v_h)_{0,T}) - \sum_{T \in \mathcal{F}_h(\Omega)} (\mathbf{n}_{\partial T} \cdot \hat{\boldsymbol{\sigma}}_h, v_h)_{0,\partial T} = 0. \quad (2.10b)$$

Here, the PWDG flux functions \hat{u}_h and $\hat{\boldsymbol{\sigma}}_h$ are given by

$$\hat{u}_h|_E := \begin{cases} \{u_h\}_E - \frac{\beta}{i\omega} [\nabla u_h]_E, & E \in \mathcal{E}_h(\Omega) \\ u_h - \delta \left(\frac{1}{i\omega} \mathbf{n}_E \cdot \nabla u_h + u_h - \frac{1}{i\omega} g \right), & E \in \mathcal{E}_h(\Gamma) \end{cases}, \quad (2.11a)$$

$$\hat{\boldsymbol{\sigma}}_h|_E := \begin{cases} \frac{1}{i\omega} \{\nabla u_h\}_E - \alpha [u_h]_E, & E \in \mathcal{E}_h(\Omega) \\ \frac{1}{i\omega} \nabla u_h - (1 - \delta) \left(\frac{1}{i\omega} \nabla u_h + \mathbf{n}_E u_h - \frac{1}{i\omega} \mathbf{n}_E g \right), & E \in \mathcal{E}_h(\Gamma) \end{cases}, \quad (2.11b)$$

where \mathbf{n}_E is the exterior unit normal on E and $\alpha > 0$, $\beta > 0$, $\delta \in (0, 1)$ are flux parameters independent of h , p , and ω .

By choosing $\boldsymbol{\tau}_h = \nabla v_h$ in (2.10a), we can eliminate $\boldsymbol{\sigma}_h$ from (2.10a), (2.10b) and obtain the following primal variational formulation of the PWDG method: Find $u_h \in V_h$ such that for all $v_h \in V_h$ it holds

$$\begin{aligned} & \sum_{T \in \mathcal{T}_h(\Omega)} ((\nabla u_h, \nabla v_h)_{0,T} - \omega^2 (u_h, v_h)_{0,T}) \\ & - \sum_{T \in \mathcal{T}_h(\Omega)} ((u_h - \hat{u}_h, \mathbf{n}_{\partial T} \cdot \nabla v_h)_{0,\partial T} + i\omega (\mathbf{n}_{\partial T} \cdot \hat{\boldsymbol{\sigma}}_h, v_h)_{0,\partial T}) = 0. \end{aligned} \quad (2.12)$$

Moreover, using Green's formula for the first term on the left-hand side in (2.12) and observing $(-\Delta - \omega^2 I)u_h|_T = 0$, $T \in \mathcal{T}_h(\Omega)$, we are led to a formulation of the PWDG method involving only integrals over edges $E \in \mathcal{E}_h(\bar{\Omega})$: Find $u_h \in V_h$ such that

$$a_h(u_h, v_h) = \ell_h(v_h), \quad v_h \in V_h, \quad (2.13)$$

where the sesquilinear form $a_h(\cdot, \cdot) : V_h \times V_h \rightarrow \mathbb{C}$ and the functional $\ell_h : V_h \rightarrow \mathbb{C}$ are given by

$$\begin{aligned} a_h(u_h, v_h) := & \sum_{E \in \mathcal{E}_h(\Omega)} ((\{u_h\}_E, \mathbf{n}_E \cdot [\nabla v_h]_E)_{0,E} + i\beta\omega^{-1} (\mathbf{n}_E \cdot [\nabla u_h]_E, \mathbf{n}_E \cdot [\nabla v_h]_E)_{0,E} \\ & - (\mathbf{n}_E \cdot \{\nabla u_h\}_E, [v_h]_E)_{0,E} + i\alpha\omega (\{u_h\}_E, [v_h]_E)_{0,E}) \\ & \sum_{E \in \mathcal{E}_h(\Gamma)} ((1 - \delta)(u_h, \mathbf{n}_E \cdot \nabla v_h)_{0,E} + i\delta\omega^{-1} (\mathbf{n}_E \cdot \nabla u_h, \mathbf{n}_E \cdot \nabla v_h)_{0,E} \\ & - \delta(\mathbf{n}_E \cdot \nabla u_h, v_h)_{0,E} + i(1 - \delta)\omega (u_h, v_h)_{0,E}), \end{aligned} \quad (2.14a)$$

$$\ell_h(v_h) := \sum_{E \in \mathcal{E}_h(\Gamma)} (i\delta\omega^{-1} (g, \mathbf{n}_E \cdot \nabla v_h)_{0,E} + (1 - \delta)(g, v_h)_{0,E}). \quad (2.14b)$$

As has been shown in [11], the variational equation (2.13) admits a unique solution $u_h \in V_h$. Moreover, if the solution u of (2.1a), (2.1b) satisfies $u \in H^{k+1}(\Omega)$, $k \in \mathbb{N}$, and if the mesh width h of the triangulation $\mathcal{T}_h(\Omega)$ satisfies $\omega h \leq \kappa$ for some $\kappa > 0$, then there exists a constant $C > 0$, independent of p and u , but depending on κ , such that the following a priori error estimate holds true (cf. Theorem 3.14 in [11])

$$\|u - u_h\|_{0,\Omega} \leq C\omega^{-1} \text{diam}(\Omega) h^{k-1} \left(\frac{\log p}{p}\right)^{k-1/2} \|u\|_{k+1,\omega,\Omega}, \quad (2.15)$$

where $\|\cdot\|_{k+1,\omega,\Omega}$ stands for the ω -weighted Sobolev norm

$$\|v\|_{k+1,\omega,\Omega} := \left(\sum_{j=0}^{k+1} \omega^{2(k+1-j)} |v|_{j,\Omega}^2 \right)^{1/2}, \quad v \in H^{k+1}(\Omega).$$

The global PWDG space V_h is spanned by Np basis functions

$$\begin{aligned}
V_h &= \text{span}(\varphi_h^{(1)}, \dots, \varphi_h^{(Np)}), \\
\varphi_h^{((j-1)p+\ell)} &:= \exp(i\omega(\cos(\theta_{(j-1)p+\ell}), \sin(\theta_{(j-1)p+\ell}))^T \cdot \mathbf{x}), \\
1 \leq j \leq N, 1 \leq \ell \leq p.
\end{aligned} \tag{2.16}$$

Then, $u_h \in V_h$ can be written as

$$u_h = \sum_{j=1}^{Np} u_j \varphi_h^{(j)}, \quad u_j \in \mathbb{C}, 1 \leq j \leq Np. \tag{2.17}$$

Further, setting $\mathbf{y} := (y_1, \dots, y_{Np})^T \in \mathbb{C}^{Np}$ with $y_j := u_j$, $1 \leq j \leq Np$, the PWDG approximation (2.13) represents a complex linear algebraic system

$$\mathbf{A}(\boldsymbol{\theta})\mathbf{y} = \mathbf{b}(\boldsymbol{\theta}), \tag{2.18}$$

where the matrix $\mathbf{A}(\boldsymbol{\theta}) = (a_{k\ell}(\boldsymbol{\theta}))_{k,\ell=1}^{Np} \in \mathbb{C}^{Np \times Np}$ and the vector $\mathbf{b}(\boldsymbol{\theta}) = (b_1(\boldsymbol{\theta}), \dots, b_{Np}(\boldsymbol{\theta}))^T \in \mathbb{C}^{Np}$ are given by

$$\begin{aligned}
a_{k\ell}(\boldsymbol{\theta}) &:= a_h(\varphi_h^{(\ell)}(\boldsymbol{\theta}), \varphi_h^{(k)}(\boldsymbol{\theta})), \quad 1 \leq k, \ell \leq Np, \\
b_\ell(\boldsymbol{\theta}) &:= \ell_h(\varphi_h^{(\ell)}), \quad 1 \leq \ell \leq Np.
\end{aligned} \tag{2.19}$$

3. Optimization of the plane wave directions

The a priori estimate (2.15) for the L^2 -norm of the global discretization error tells us how the error depends on the number p of plane wave directions, but it does not provide any information on the appropriate choice of the directions $\mathbf{d}_{(j-1)p+\ell} = (\cos(\theta_{(j-1)p+\ell}), \sin(\theta_{(j-1)p+\ell}))^T$, $1 \leq j \leq N$, $1 \leq \ell \leq p$, except that they are supposed to satisfy assumption (2.8). In fact, since

$$V_h = \text{span}(\exp(i\omega \mathbf{d}_1 \cdot \mathbf{x}), \dots, \exp(i\omega \mathbf{d}_{Np} \cdot \mathbf{x})), \tag{3.1}$$

the solution $u_h \in V_h$ of (2.13) depends on $\boldsymbol{\theta} := (\theta_1, \dots, \theta_{Np})^T \in \mathbf{K}$ according to

$$u_h(\boldsymbol{\theta}) = \sum_{k=1}^N \sum_{\ell=1}^p u_{k\ell} \exp(i\omega \mathbf{d}_{(k-1)p+\ell} \cdot \mathbf{x})|_{T_k}, \quad u_{k\ell} \in \mathbb{C}. \tag{3.2}$$

We attempt to choose $\boldsymbol{\theta} \in \mathbf{K}$ such that a given continuously differentiable objective functional $J : V_h \times \mathbb{C}^{Np} \rightarrow \mathbb{R}$ is minimized. This can be formulated as the optimal control problem

$$\min_{u_h \in V_h, \boldsymbol{\theta} \in \mathbf{K}} J(u_h, \boldsymbol{\theta}), \tag{3.3a}$$

subject to the PWDG constraint

$$a_h(u_h(\theta), v_h(\theta)) = \ell_h(v_h(\theta)), \quad v_h(\theta) \in V_h. \quad (3.3b)$$

The algebraic formulation of (3.3a), (3.3b) turns out to be

$$\min_{\mathbf{y} \in \mathbb{C}^{Np}, \boldsymbol{\theta} \in \mathbf{K}} J(\mathbf{y}, \boldsymbol{\theta}), \quad (3.4a)$$

subject to the state equation

$$e(\mathbf{y}, \boldsymbol{\theta}) := \mathbf{A}(\boldsymbol{\theta})\mathbf{y} - \mathbf{b}(\boldsymbol{\theta}) = 0. \quad (3.4b)$$

Remark 3.1. A particular choice of the objective functional J is

$$J(u_h, \boldsymbol{\theta}) := \frac{1}{2} \|u_h(\boldsymbol{\theta}) - u^d\|_{0,\Omega}^2, \quad (3.5)$$

where $u^d \in L^2(\Omega)$ is a given objective functional. Introducing the Hermitian matrix $\mathbf{M}(\boldsymbol{\theta}) = (m_{k\ell}(\boldsymbol{\theta}))_{k,\ell=1}^{Np} \in \mathbb{C}^{Np \times Np}$ and the vector $\mathbf{c}(\boldsymbol{\theta}) = (c_1(\boldsymbol{\theta}), \dots, c_{Np}(\boldsymbol{\theta}))^T$ according to

$$\begin{aligned} m_{k\ell}(\boldsymbol{\theta}) &:= (\varphi_h^{(k)}, \varphi_h^{(\ell)})_{0,\Omega}, \quad 1 \leq k, \ell \leq Np, \\ c_\ell(\boldsymbol{\theta}) &:= (u^d, \varphi_h^{(\ell)})_{0,\Omega}, \quad 1 \leq \ell \leq Np, \end{aligned} \quad (3.6)$$

in algebraic form the objective functional reads as follows

$$J(\mathbf{y}, \boldsymbol{\theta}) := \frac{1}{2} \langle \mathbf{M}(\boldsymbol{\theta})\mathbf{y}, \mathbf{y} \rangle - \text{Re}(\langle \mathbf{c}(\boldsymbol{\theta}), \mathbf{y} \rangle). \quad (3.7)$$

We denote by $\mathbf{G} : \mathbf{K} \rightarrow \mathbb{C}^{Np}$ the control-to-state map which assigns to the control $\boldsymbol{\theta} \in \mathbf{K}$ the unique solution $\mathbf{y} \in \mathbb{C}^{Np}$ of the state equation (3.4b) and by $J_{\text{red}} : \mathbf{K} \rightarrow \mathbb{R}$ the reduced objective functional

$$J_{\text{red}}(\boldsymbol{\theta}) := J(\mathbf{G}(\boldsymbol{\theta}), \boldsymbol{\theta}).$$

Then, the control-reduced formulation of the optimal control problem (3.4a), (3.4b) reads as follows

$$\min_{\boldsymbol{\theta} \in \mathbf{K}} J_{\text{red}}(\boldsymbol{\theta}). \quad (3.8)$$

The existence of a solution follows by standard arguments from the theory of constrained finite dimensional optimization problems (cf., e.g., [16]).

Theorem 3.2. *The optimal control problem (3.4a), (3.4b) admits an optimal solution $(\mathbf{y}^*, \boldsymbol{\theta}^*) \in \mathbb{C}^{Np} \times \mathbf{K}$.*

Proof. Let $\{\boldsymbol{\theta}^{(n)}\}_{n \in \mathbb{N}}$, $\boldsymbol{\theta}^{(n)} \in \mathbf{K}$, $n \in \mathbb{N}$, be a minimizing sequence, i.e., it holds

$$J_{\text{red}}(\boldsymbol{\theta}^{(n)}) \rightarrow \min_{\boldsymbol{\theta} \in \mathbf{K}} J_{\text{red}}(\boldsymbol{\theta}) \quad \text{as } n \rightarrow \infty. \quad (3.9)$$

Obviously, the sequence $\{\boldsymbol{\theta}^{(n)}\}_{n \in \mathbb{N}}$ is bounded and hence, there exist a subsequence $\mathbb{N}' \subset \mathbb{N}$ and $\boldsymbol{\theta}^* \in \mathbb{R}^p$ such that

$$\boldsymbol{\theta}^{(n)} \rightarrow \boldsymbol{\theta}^*, \quad \mathbb{N}' \ni n \rightarrow \infty.$$

In view of the closedness of \mathbf{K} , we have $\boldsymbol{\theta}^* \in \mathbf{K}$. Moreover, due to the continuity of both the control-to-state map \mathbf{G} and of the reduced objective functional J_{red} we deduce

$$\mathbf{G}(\boldsymbol{\theta}^{(n)}) \rightarrow \mathbf{G}(\boldsymbol{\theta}^*), \quad J_{\text{red}}(\boldsymbol{\theta}^{(n)}) \rightarrow J_{\text{red}}(\boldsymbol{\theta}^*) \quad \mathbb{N}' \ni n \rightarrow \infty.$$

Consequently, from (3.9) we have

$$J_{\text{red}}(\boldsymbol{\theta}^*) = \min_{\boldsymbol{\theta} \in \mathbf{K}} J_{\text{red}}(\boldsymbol{\theta}),$$

and with $\mathbf{y}^* := \mathbf{G}(\boldsymbol{\theta}^*)$ it follows that the pair $(\mathbf{y}^*, \boldsymbol{\theta}^*) \in \mathbb{C}^{Np} \times \mathbf{K}$ is an optimal solution of (3.4a), (3.4b). \square

Remark 3.3. Since the control-to-state map \mathbf{G} is a non-convex function of the control $\boldsymbol{\theta}$, we do not have uniqueness of an optimal solution.

4. First order necessary optimality conditions

We will derive the first order necessary optimality conditions for the optimal control problem (3.4a), (3.4b) by the method of Lagrange multipliers which is justified if the linear independence constraint qualification holds true [14], [16]. To this end, we note that the bound constraints on the control can be expressed as the inequalities $\mathbf{g}(\boldsymbol{\theta}) \leq \mathbf{0}$, where the mapping $\mathbf{g} = (\mathbf{g}_1, \mathbf{g}_2) : \mathbb{R}^{Np} \rightarrow \mathbb{R}^{Np} \times \mathbb{R}^{Np}$ is defined by means of

$$\begin{aligned} \mathbf{g}_1(\boldsymbol{\theta}) &:= (\hat{\theta}_2 - \hat{\theta}_1 - \theta_{\max}, \dots, \hat{\theta}_{Np+1} - \hat{\theta}_{Np} - \theta_{\max}), \\ \mathbf{g}_2(\boldsymbol{\theta}) &:= (\theta_{\min} - (\hat{\theta}_2 - \hat{\theta}_1), \dots, \theta_{\min} - (\hat{\theta}_{Np+1} - \hat{\theta}_{Np})). \end{aligned} \quad (4.1)$$

For a local minimum $(\mathbf{y}^*, \boldsymbol{\theta}^*) \in \mathbb{C}^{Np} \times \mathbf{K}$ of (3.4a), (3.4b) the active set is given by $A(\boldsymbol{\theta}^*) = A_1(\boldsymbol{\theta}^*) \cup A_2(\boldsymbol{\theta}^*)$ where

$$A_1(\boldsymbol{\theta}^*) := \{q \in \{1, \dots, p\} \mid \hat{\theta}_{q+1}^* - \hat{\theta}_q^* - \theta_{\max} = 0\}, \quad (4.2a)$$

$$A_2(\boldsymbol{\theta}^*) := \{q \in \{1, \dots, p\} \mid \theta_{\min} - (\hat{\theta}_{q+1}^* - \hat{\theta}_q^*) = 0\}, \quad (4.2b)$$

where $\hat{\theta}_q^*$ is defined as in (2.7) with θ_q replaced by θ_q^* .

We refer to $I(\boldsymbol{\theta}^*) := \{1, \dots, p\} \setminus A(\boldsymbol{\theta}^*)$ as the inactive set. The linear independence constraint qualification requires the linearization of $(e, (\mathbf{g}_1)_{A_1(\boldsymbol{\theta}^*)}, (\mathbf{g}_2)_{A_2(\boldsymbol{\theta}^*)})$ at $(\mathbf{y}^*, \boldsymbol{\theta}^*)$ to be surjective.

Theorem 4.1. *Let $p_i^* := \text{card}(A_i(\boldsymbol{\theta}^*))$, $1 \leq i \leq 2$ and assume $I(\boldsymbol{\theta}^*) \neq \emptyset$. The mapping*

$$(\nabla e(\mathbf{y}^*, \boldsymbol{\theta}^*), \nabla \mathbf{g}_{1, A_1(\boldsymbol{\theta}^*)}(\boldsymbol{\theta}^*), \nabla \mathbf{g}_{2, A_2(\boldsymbol{\theta}^*)}(\boldsymbol{\theta}^*)) : \mathbb{C}^{Np} \times \mathbb{R}^{Np} \rightarrow \mathbb{C}^{Np} \times \mathbb{R}^{p_1^*} \times \mathbb{R}^{p_2^*}$$

is surjective. In particular, for any $(\mathbf{r}, \mathbf{s}_1, \mathbf{s}_2) \in \mathbb{C}^{Np} \times \mathbb{R}^{p_1^} \times \mathbb{R}^{p_2^*}$ there exists a unique solution $(\delta \mathbf{y}, \delta \boldsymbol{\theta}) \in \mathbb{C}^{Np} \times \mathbb{R}^{Np}$ of the equation*

$$(\nabla e(\mathbf{y}^*, \boldsymbol{\theta}^*)(\delta \mathbf{y}, \delta \boldsymbol{\theta}), \nabla \mathbf{g}_{1, A_1(\boldsymbol{\theta}^*)}(\boldsymbol{\theta}^*)\delta \boldsymbol{\theta}, \nabla \mathbf{g}_{2, A_2(\boldsymbol{\theta}^*)}(\boldsymbol{\theta}^*)\delta \boldsymbol{\theta}) = (\mathbf{r}, \mathbf{s}_1, \mathbf{s}_2).$$

Proof. For $k \in A_1(\boldsymbol{\theta}^*)$ we obviously have

$$\nabla g_{1, k'}(\boldsymbol{\theta}^*) = \begin{cases} -1, & k' = k \\ +1, & k' = k + 1, \\ 0, & \text{otherwise} \end{cases} \quad (4.3)$$

whereas for $k \in A_2(\boldsymbol{\theta}^*)$

$$\nabla g_{2, k'}(\boldsymbol{\theta}^*) = \begin{cases} +1, & k' = k \\ -1, & k' = k + 1. \\ 0, & \text{otherwise} \end{cases} \quad (4.4)$$

Since $I(\boldsymbol{\theta}^*) \neq \emptyset$, there exists $q \in \{1, \dots, Np\}$ such that $q \in I(\boldsymbol{\theta}^*)$. We renumber the controls according to $\tilde{\theta}_k^* := \theta_{q+k-1}^*$, $\tilde{\theta}_{k+p}^* = \hat{\theta}_k^* + 2\pi$, $1 \leq k \leq Np$, and set $(\delta \boldsymbol{\theta})_k = 0$ for $k \in I(\tilde{\boldsymbol{\theta}}^*)$. If $A(\tilde{\boldsymbol{\theta}}^*) = \emptyset$, there is nothing to show. If $A(\tilde{\boldsymbol{\theta}}^*) \neq \emptyset$, there exists

$$k_{\min} := \min\{k \in \{2, \dots, Np\} \mid k \in A(\tilde{\boldsymbol{\theta}}^*)\}.$$

Moreover, in view of $Np + 1 \in I(\tilde{\boldsymbol{\theta}}^*)$, there also exists

$$k_{\max} := \min\{k \in \{k_{\min} + 1, \dots, Np + 1\} \mid k \in I(\tilde{\boldsymbol{\theta}}^*)\}.$$

In view of (4.3), (4.4), $(\delta\theta)_k$, $k_{\min} \leq k \leq k_{\max} - 1$, is the unique solution of a linear algebraic system with a regular upper triangular matrix. For the computation of $(\delta\theta)_k \in A(\theta^*) \setminus \{k_{\min}, \dots, k_{\max} - 1\}$ we proceed in the same way.

On the other hand, the equation $\nabla e(\mathbf{y}^*, \theta^*)(\delta\mathbf{y}, \delta\theta) = \mathbf{r}$ can be equivalently written as

$$\mathbf{A}(\theta)\delta\mathbf{y} = \nabla_{\theta}(\mathbf{b}(\theta^*) - \mathbf{A}(\theta^*)\mathbf{y}^*)\delta\theta,$$

which has a unique solution $\delta\mathbf{y} \in \mathbb{C}^{Np}$. □

Due to Theorem 4.1, the necessary optimality conditions can be derived by the method of Lagrange multipliers.

Theorem 4.2. *Assume that $(\mathbf{y}^*, \theta^*) \in \mathbb{C}^{Np} \times \mathbf{K}$ is an optimal solution of (3.4a), (3.4b). Then there exist an adjoint state $\mathbf{p}^* \in \mathbb{C}^{Np}$ and a multiplier $\boldsymbol{\mu}^* = (\boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*) \in \mathbb{R}_+^{2Np}$, $\boldsymbol{\mu}_i^* = (\mu_{i,1}^*, \dots, \mu_{i,Np}^*)^T$, $1 \leq i \leq 2$, such that the state equation, the adjoint state equation and the gradient equation*

$$\mathbf{A}(\theta^*)\mathbf{y}^* - \mathbf{b}(\theta^*) = \mathbf{0},$$

$$\mathbf{A}^H(\theta^*)\mathbf{p}^* + \mathbf{J}_{\mathbf{y}}(\mathbf{y}^*, \theta^*) = \mathbf{0},$$

$$\nabla_{\theta}J(\mathbf{y}^*, \theta^*) + \text{Re}(\langle \nabla_{\theta}(\mathbf{A}(\theta^*)\mathbf{y}^* - \mathbf{b}(\theta^*)), \mathbf{p}^* \rangle) + \nabla_{\theta}\mathbf{g}_1(\theta^*)^T \boldsymbol{\mu}_1^* + \nabla_{\theta}\mathbf{g}_2(\theta^*)^T \boldsymbol{\mu}_2^* = \mathbf{0}$$

are satisfied as well as the complementarity conditions

$$g_{i,q}(\theta^*) \leq 0, \mu_{i,q}^* \geq 0, g_{i,q}(\theta^*)\mu_{i,q}^* = 0, \quad 1 \leq q \leq Np, 1 \leq i \leq 2.$$

Proof. We introduce the Lagrangian $L : \mathbb{C}^{Np} \times \mathbb{R}^{Np} \times \mathbb{C}^{Np} \times \mathbb{R}_+^{2Np}$ according to

$$L(\mathbf{y}, \theta, \mathbf{p}, \boldsymbol{\mu}) := J(\mathbf{y}, \theta) + \text{Re}(\langle e(\mathbf{y}, \theta), \mathbf{p} \rangle) + \mathbf{g}_1(\theta)^T \boldsymbol{\mu}_1 + \mathbf{g}_2(\theta)^T \boldsymbol{\mu}_2.$$

Setting $\mathbf{x} := (\mathbf{y}, \theta, \mathbf{p})$ and $\mathbf{x}^* := (\mathbf{y}^*, \theta^*, \mathbf{p}^*)$, the first order necessary optimality conditions are given by

$$\frac{\partial L}{\partial \mathbf{y}}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \mathbf{0}, \quad \frac{\partial L}{\partial \theta}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \mathbf{0}, \quad \frac{\partial L}{\partial \mathbf{p}}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \mathbf{0}, \quad (4.5a)$$

$$\frac{\partial L}{\partial \boldsymbol{\mu}_i}(\mathbf{x}^*, \boldsymbol{\mu}^*)^T (\mathbf{v}_i - \boldsymbol{\mu}_i^*) \leq 0, \quad \mathbf{v}_i \in \mathbb{R}_+^{Np}, 1 \leq i \leq 2. \quad (4.5b)$$

The state equation, the adjoint state equation, and the gradient equation result from the third, first, and second equation in (4.5a), whereas the complementarity conditions are a consequence of (4.5b). □

5. Projected gradient method

The projected gradient method is based on the formulation of the gradient equation as the variational inequality

$$-\nabla_{\theta} J(\mathbf{y}^*, \boldsymbol{\theta}^*) + \text{Re}(\langle \nabla_{\theta}(\mathbf{b}(\boldsymbol{\theta}^*) - \mathbf{A}(\boldsymbol{\theta}^*)\mathbf{y}^*), \mathbf{p}^* \rangle) \in \partial I_{\mathbf{K}},$$

where $\partial I_{\mathbf{K}}$ is the subdifferential of the indicator function of the constraint set \mathbf{K} .

Projected gradient method

Step 1: Choose an initial control $\boldsymbol{\theta}^{(0)} \in \mathbf{K}$ and a tolerance $TOL > 0$ and set $n = 0$.

Step 2.1: Set $n = n + 1$ and compute $\mathbf{y}^{(n)} \in \mathbb{C}^{Np}$ and $\mathbf{p}^{(n)} \in \mathbb{C}^{Np}$ as the unique solutions of the state equation

$$\mathbf{A}(\boldsymbol{\theta}^{(n-1)})\mathbf{y}^{(n)} = \mathbf{b}(\boldsymbol{\theta}^{(n-1)})$$

and of the adjoint state equation

$$\mathbf{A}^H(\boldsymbol{\theta}^{(n-1)})\mathbf{p}^{(n)} = -J_{\mathbf{y}}(\mathbf{y}^{(n)}, \boldsymbol{\theta}^{(n-1)}).$$

Step 2.2: Compute $\tilde{\boldsymbol{\theta}}^{(n)} \in \mathbb{R}^p$ according to

$$\tilde{\boldsymbol{\theta}}^{(n)} = \boldsymbol{\theta}^{(n-1)} - \kappa(\nabla_{\theta} J(\mathbf{y}^{(n)}, \boldsymbol{\theta}^{(n-1)}) + \text{Re}(\langle \nabla_{\theta}(\mathbf{A}(\boldsymbol{\theta}^{(n-1)})\mathbf{y}^{(n)} - \mathbf{b}(\boldsymbol{\theta}^{(n-1)})), \mathbf{p}^{(n)} \rangle)),$$

where $\kappa > 0$ is the Armijo line search parameter.

Step 2.3: Compute $\boldsymbol{\theta}^{(n)}$ as the projection of $\tilde{\boldsymbol{\theta}}^{(n)}$ onto the constraint set \mathbf{K} .

Step 2.4: If $n > 1$ and

$$|J(\mathbf{y}^{(n)}, \boldsymbol{\theta}^{(n)}) - J(\mathbf{y}^{(n-1)}, \boldsymbol{\theta}^{(n-1)})| < TOL,$$

stop the algorithm. Otherwise, go to Step 2.1.

We will provide some details regarding the numerical realization of **Step 2.2** in case the objective functional J is chosen as in (3.5). For the update formula we need to compute the following quantity:

$$\nabla_{\theta} J(\mathbf{y}, \boldsymbol{\theta}) + \text{Re}(\langle \nabla_{\theta} \mathbf{A}(\boldsymbol{\theta})\mathbf{y} - \mathbf{b}(\boldsymbol{\theta}), \mathbf{p} \rangle).$$

For the computation of $\nabla_{\theta} J(\mathbf{y}, \boldsymbol{\theta})$ we observe (3.7). Moreover, according to (3.4b) the vector $\mathbf{y} = (y_1, \dots, y_{Np})^T$ is the unique solution of

$$\mathbf{A}(\boldsymbol{\theta})\mathbf{y} = \mathbf{b}(\boldsymbol{\theta}),$$

where the $Np \times Np$ matrix $A(\boldsymbol{\theta})$ and the Np vector $\mathbf{b}(\boldsymbol{\theta})$ are given by (2.19). We note that for any two given basis functions $\phi_h^{(k)}$ and $\phi_h^{(\ell)}$ either,

$$\mu(\text{supp}(\phi_h^{(k)}) \cap \text{supp}(\phi_h^{(\ell)})) = 0$$

or,

$$\text{supp}(\phi_h^{(k)}) \cap \text{supp}(\phi_h^{(\ell)}) = T \in \mathcal{T}_h(\Omega),$$

where μ is the 2-D Lebesgue measure. Let $T_{k,\ell}$, $1 \leq k, \ell \leq Np$, be defined as

$$T_{k,\ell} := \begin{cases} \emptyset, & \text{if } \mu(\text{supp}(\phi_h^{(k)}) \cap \text{supp}(\phi_h^{(\ell)})) = 0, \\ \text{supp}(\phi_h^{(k)}) \cap \text{supp}(\phi_h^{(\ell)}), & \text{otherwise} \end{cases},$$

and let T_ℓ , $1 \leq \ell \leq Np$ be given by

$$T_\ell := \text{supp}(\phi_h^{(\ell)}) \in \mathcal{T}_h(\Omega).$$

Hence, we can rewrite (3.6) as

$$\begin{aligned} m_{k\ell}(\boldsymbol{\theta}) &:= \int_{T_{k,\ell}} \exp(i\omega \mathbf{d}_k \cdot \mathbf{x}) \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x}, & 1 \leq k, \ell \leq Np, \\ c_\ell(\boldsymbol{\theta}) &:= \int_{T_\ell} u_d \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x}, & 1 \leq \ell \leq Np. \end{aligned} \quad (5.1)$$

In view of (??) we obtain

$$\nabla_\theta J(\mathbf{y}, \boldsymbol{\theta}) = \nabla_\theta \left(\frac{1}{2} \sum_{k,\ell=1}^{Np} m_{k\ell}(\boldsymbol{\theta}) y_k \bar{y}_\ell \right) - \nabla_\theta \left(\text{Re} \sum_{k=1}^{Np} c_k(\boldsymbol{\theta}) \bar{y}_k \right). \quad (5.2)$$

Differentiating (5.1) with respect to θ_j it follows that

$$\begin{aligned} & \frac{\partial}{\partial \theta_j} \left(\frac{1}{2} \sum_{k,\ell=1}^{Np} m_{k\ell}(\boldsymbol{\theta}) y_\ell \bar{y}_k \right) \\ &= \frac{1}{2} \left(\sum_{\ell=1}^{Np} y_\ell \bar{y}_\ell \int_{T_{j,\ell}} (i\omega \mathbf{d}_j^* \cdot \mathbf{x}) \exp(i\omega \mathbf{d}_j \cdot \mathbf{x}) \cdot \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x} \right. \\ & \quad \left. + \sum_{\ell=1}^{Np} y_\ell \bar{y}_j \int_{T_{j,\ell}} (-i\omega \mathbf{d}_j^* \cdot \mathbf{x}) \exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x}) \cdot \overline{\exp(i\omega \mathbf{d}_j \cdot \mathbf{x})} d\mathbf{x} \right) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \left(\sum_{\ell=1}^{Np} y_j \bar{y}_\ell \int_{T_{j,\ell}} (i\omega \mathbf{d}_j^* \cdot \mathbf{x}) \exp(i\omega \mathbf{d}_j \cdot \mathbf{x}) \cdot \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x} \right. \\
&\quad \left. + \sum_{\ell=1}^{Np} y_j \bar{y}_\ell \int_{T_{j,\ell}} (i\omega \mathbf{d}_j^* \cdot \mathbf{x}) \exp(i\omega \mathbf{d}_j \cdot \mathbf{x}) \cdot \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x} \right) \\
&= \operatorname{Re} \sum_{\ell=1}^{Np} y_j \bar{y}_\ell \int_{T_{j,\ell}} (i\omega \mathbf{d}_j^* \cdot \mathbf{x}) \exp(i\omega \mathbf{d}_j \cdot \mathbf{x}) \cdot \overline{\exp(i\omega \mathbf{d}_\ell \cdot \mathbf{x})} d\mathbf{x}, \tag{5.3}
\end{aligned}$$

where $\mathbf{d}_j^* = (-\sin(\theta_j), \cos(\theta_j))^T$. Moreover, we obtain

$$\frac{\partial}{\partial \theta_j} \left(\operatorname{Re} \sum_{k=1}^{Np} c_k(\boldsymbol{\theta}) \bar{y}_k \right) = \operatorname{Re} \left(\bar{y}_j \int_{\Omega} (-i\omega \mathbf{d}_j^* \cdot \mathbf{x}) u_d \overline{\exp(i\omega \mathbf{d}_j \cdot \mathbf{x})} d\mathbf{x} \right). \tag{5.4}$$

On the other hand, for $\operatorname{Re}(\nabla_{\boldsymbol{\theta}} \langle \mathbf{A}(\boldsymbol{\theta}) \mathbf{y} - \mathbf{b}(\boldsymbol{\theta}), \mathbf{p} \rangle)$ we have

$$\begin{aligned}
\operatorname{Re} \left(\frac{\partial}{\partial \theta_j} \langle \mathbf{A}(\boldsymbol{\theta}) \mathbf{y} - \mathbf{b}(\boldsymbol{\theta}), \mathbf{p} \rangle \right) &= \operatorname{Re} \left(\frac{\partial}{\partial \theta_j} \sum_{k,\ell=1}^{Np} (a_{k\ell}(\boldsymbol{\theta}) y_\ell - b_k(\boldsymbol{\theta})) \bar{p}_k \right) \\
&= \operatorname{Re} \left(\sum_{k,\ell=1}^{Np} \left(\frac{\partial a_{k\ell}(\boldsymbol{\theta})}{\partial \theta_j} y_\ell - \frac{\partial b_k(\boldsymbol{\theta})}{\partial \theta_j} \right) \bar{p}_k \right). \tag{5.5}
\end{aligned}$$

We obtain the derivatives $\frac{\partial a_{kl}(\boldsymbol{\theta})}{\partial \theta_j}$ and $\frac{\partial b_k(\boldsymbol{\theta})}{\partial \theta_j}$ by directly differentiating the formulas in (2.19).

Using (5.3)–(5.5) provides the update formula in **Step 2.2** of the projected gradient method.

6. Numerical results

We provide a documentation of numerical results illustrating the performance of the optimization algorithm by studying the global discretization error $u - u_h$ in the L^2 -norm. Therefore, we consider two examples where either the exact solution or a sufficiently accurate approximate solution is known and we choose the objective functional (3.5) with u^d being given by the exact (approximate) solution. We note that for problems where a desired state u^d is not easily available we may use any continuously differentiable objective functional (e.g., the acoustic energy associated with the problem) as outlined in Sections 3 and 4.

Example 1. As in [10], [11] we consider the Helmholtz problem (2.1a), (2.1b) in $\Omega = (0, 1) \times (-0.5, +0.5)$ with $\omega = 10$ and g in (2.1b) being chosen such that the

exact solution u (in polar coordinates) is given by

$$u(r, \varphi) = J_\xi(\omega r) \cos(\xi\varphi), \quad \xi \geq 0,$$

where J_ξ stands for the Bessel function of the first kind and order ξ . We note that for $\xi \in \mathbb{N}$ the solution is regular, whereas for $\xi \notin \mathbb{N}$ the solution satisfies $u \in H^{1+\xi-\varepsilon}(\Omega)$ for any $\varepsilon > 0$ and its derivatives have a singularity at the origin. Figure 1 displays the exact solution for $\xi = 1$ (top right), $\xi = 2/3$ (bottom left), and $\xi = 3/2$ (bottom right).

For $\xi = 1$, $\xi = 2/3$, and $\xi = 3/2$ the PWDG method has been implemented with respect to a geometrically conforming simplicial triangulation $\mathcal{T}_h(\Omega)$ consisting of eight isosceles triangles (cf. Figure 1 (top left)). The parameters α , β , and δ in the PWDG method (2.10a), (2.10b) are chosen either according to

$$\alpha = \beta = \delta = 0.5 \tag{6.1}$$

as in the ultraweak variational formulation by Cessenat and Després [4] or by means of

$$\alpha = \beta^{-1} = \delta^{-1} = \frac{10p}{\omega h \log(p)} \tag{6.2}$$

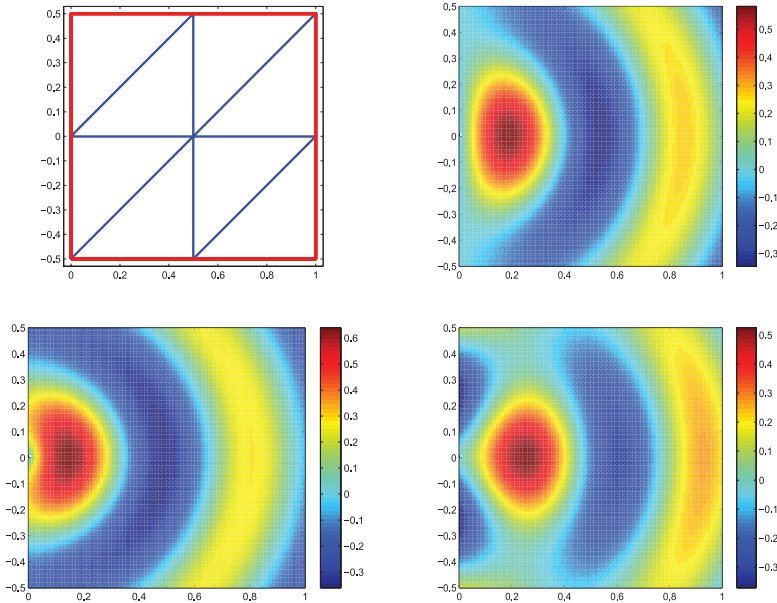


Figure 1. Example 1: The computational domain Ω and the simplicial triangulation $\mathcal{T}_h(\Omega)$ for the PWDG method (top left). The exact solution for $\xi = 1$ (top right), $\xi = 2/3$ (bottom left), and $\xi = 3/2$ (bottom right).

as suggested in [11]. For the optimization of the plane wave directions, the objective functional has been chosen as in (3.5) with the desired state u^d given by the exact solution. Moreover, we have chosen equidistributed directions with respect to the triangles $T \in \mathcal{T}_h(\Omega)$, i.e., $\mathbf{d}_{(j-1)p+\ell} = \mathbf{d}_\ell$, $1 \leq \ell \leq p$, for all $j \in \{1, \dots, N\}$. For the projected gradient method the initial distribution θ_0 has been chosen as either uniform as in [11] or random as suggested by Cessenat and Després. The state equation as well as the adjoint state equation in Step 2.1 have been solved by Gaussian elimination. We note that the optimization problem has multiple local minima and hence, starting at different initial distributions the algorithm may terminate at different local minima.

For $\xi = 1$ Figure 2 and for $\xi = 2/3$ Figure 3 display the global discretization error $u - u_h$ in the L^2 -norm $\|\cdot\|_{0,\Omega}$ as a function of the number p of plane wave basis functions. For the choice of the parameters α , β , and δ according to (6.1), the results for a uniform initial distribution θ_0 of the plane wave directions are top left and for a random initial distribution θ_0 they are shown top right. On the other hand, if the parameters α , β , and δ are chosen by means of (6.2), the results

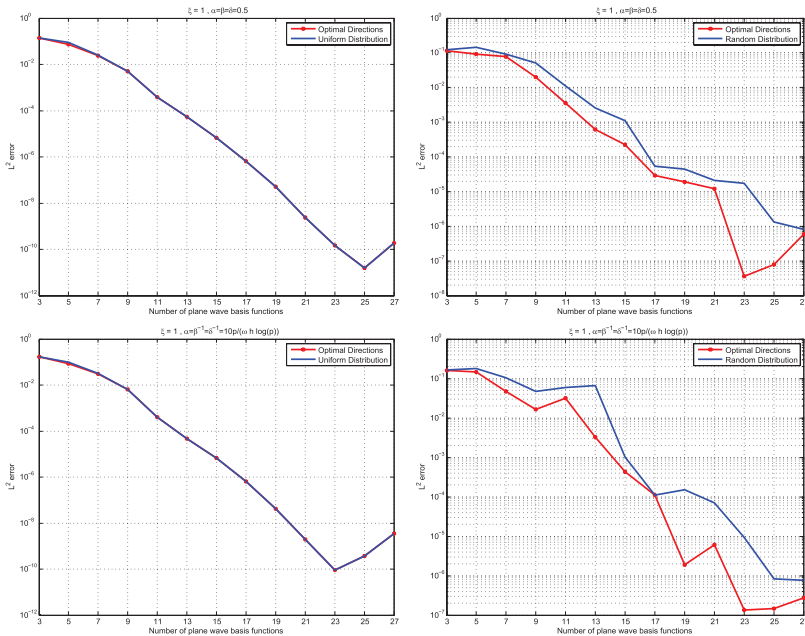


Figure 2. Example 1: The L^2 -error $\|u - u_h\|_{0,\Omega}$ as a function of the number p of plane wave basis functions for $\xi = 1$: Parameter choice (6.1) and uniform initial distribution θ_0 (top left), parameter choice (6.1) and random initial distribution θ_0 (top right), parameter choice (6.2) and uniform initial distribution θ_0 (bottom left), parameter choice (6.2) and random initial distribution θ_0 (bottom right).

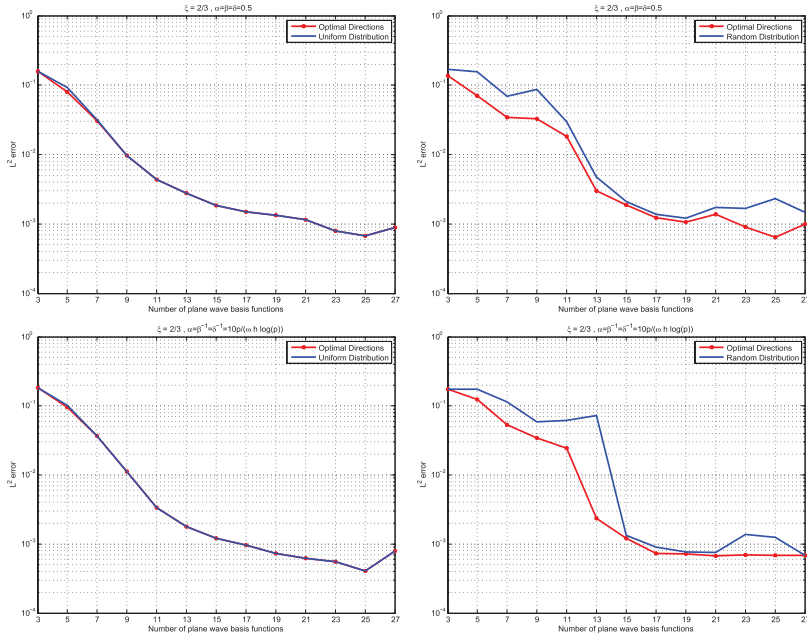


Figure 3. Example 1: The L^2 -error $\|u - u_h\|_{0,\Omega}$ as a function of the number p of plane wave basis functions for $\xi = 2/3$: Parameter choice (6.1) and uniform initial distribution θ_0 (top left), parameter choice (6.1) and random initial distribution θ_0 (top right), parameter choice (6.2) and uniform initial distribution θ_0 (bottom left), parameter choice (6.2) and random initial distribution θ_0 (bottom right).

for a uniform initial distribution θ_0 are displayed bottom left, whereas for a random initial distribution θ_0 they are shown bottom right. We see that both in case of a regular solution ($\xi = 1$) and of a singular solution ($\xi = 2/3$) the uniform distribution of θ_0 is optimal except for $p = 3, 5, 7, 9$ where it is almost optimal (cf. Figure 4 for $\xi = 2/3$ and $p = 5$). However, for a random initial distribution θ_0 the computed optimal distribution yields a reduction in the L^2 -error $\|u - u_h\|_{0,\Omega}$ up to one order of magnitude. Figure 5 shows the randomly chosen initial distribution and the computed optimal distribution for $\xi = 1$ and $p = 7$.

The (almost) optimality of the uniform distribution of the plane wave directions is probably due to the fact that the solution is symmetric (with respect to the x_1 -axis). Moreover, we see that the difference between the two parameter choices (6.1) and (6.2) is only marginal. The results for the case $\xi = 3/2$ are very similar and are thus omitted.

We note that the condition number of the matrix $\mathbf{A}(\theta)$ deteriorates with increasing number p of plane wave basis functions so that roundoff errors may effect

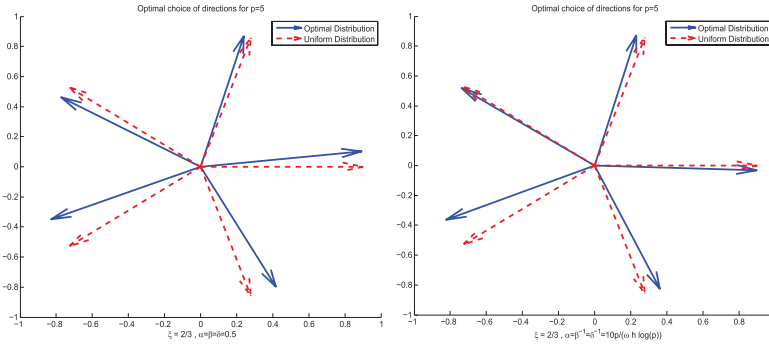


Figure 4. Example 1: Uniform initial distribution θ_0 of the plane wave directions (red dotted line) and the computed optimal distribution (blue solid line) for $\zeta = 2/3$ and $p = 5$: Parameter choice (6.1) left and parameter choice (6.2) right.

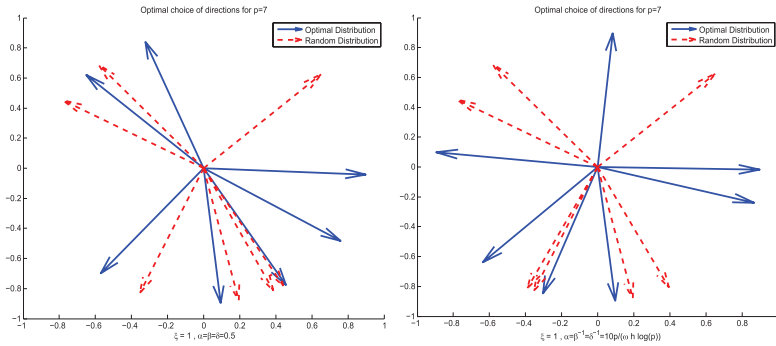


Figure 5. Example 1: Randomly chosen initial distribution θ_0 of the plane wave directions (red dotted line) and the computed optimal distribution (blue solid line) for $\zeta = 1$ and $p = 7$: Parameter choice (6.1) left and parameter choice (6.2) right.

the convergence. For $\zeta = 1$ we observe such a behavior for $p \geq 23$ (cf. Figure 2), whereas in the singular case $\zeta = 2/3$ a slowdown of the convergence can already be seen for $p \geq 17$ (cf. Figure 3).

Example 2. The second example deals with a screen problem which describes an acoustic wave scattered at a sound-soft scatterer:

$$-\Delta u - \omega^2 u = f \quad \text{in } \Omega, \tag{6.3a}$$

$$\mathbf{n} \cdot \nabla u + i\omega u = g \quad \text{on } \Gamma_R, \tag{6.3b}$$

$$u = 0 \quad \text{on } \Gamma_D, \tag{6.3c}$$

The computational domain is given by $\Omega := (-1, +1)^2 \setminus (S_1 \cup S_2)$ where

$$\begin{aligned} S_1 &:= \text{conv}((0, 0), (-0.25, +0.50), (-0.50, +0.50)), \\ S_2 &:= \text{conv}((0, 0), (+0.25, -0.50), (+0.50, -0.50)). \end{aligned}$$

Moreover, $\Gamma_R = \partial(-1, +1)^2$ and $\Gamma_D := \partial S_1 \cup \partial S_2$. The right-hand sides f and g are chosen according to $f \equiv 0$ and

$$g = \cos(\omega x_2) + i \sin(\omega x_2).$$

The exact solution u is not known explicitly. As a substitute for the exact solution we have used an approximate solution u_s computed by the adaptive Interior Penalty Discontinuous Galerkin method from [13] with a sufficiently large number of refinement steps. For $\omega = 15$, the approximate solution u_s is displayed in Figure 6 (right).

The PWDG method has been implemented with respect to a geometrically conforming simplicial triangulation $\mathcal{T}_h(\Omega)$ shown in Figure 6 (left). The parameters α , β , and δ of the PWDG method have been chosen according to (6.1) and (6.2). For the optimization, we have chosen u^d in the objective functional as the substitute solution u_s . Moreover, for the projected gradient method the initial distribution θ_0 of the plane wave directions has been chosen as a uniform distribution. In Step 2.1 of the projected gradient method, the state equation and the adjoint state equation have been solved numerically by GMRES [17].

In this example, we have compared uniform directions with optimal directions in case of the same directions for each element of the triangulation (optimal directions I) and different directions per element (optimal directions II).

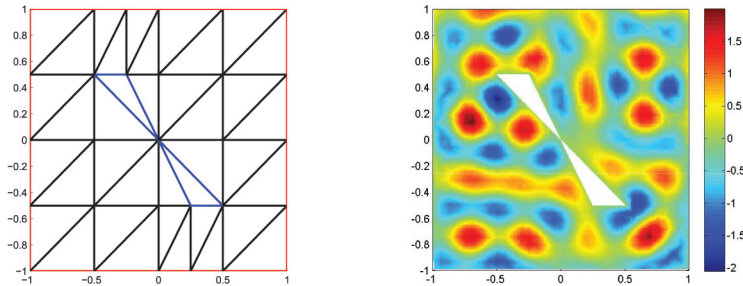


Figure 6. Example 2: The computational domain Ω and the simplicial triangulation $\mathcal{T}_h(\Omega)$ for the PWDG method (left; the sound-soft scatterer is shown in blue). The substitute solution u_s computed by the adaptive Interior Penalty Discontinuous Galerkin method from [13] (right).

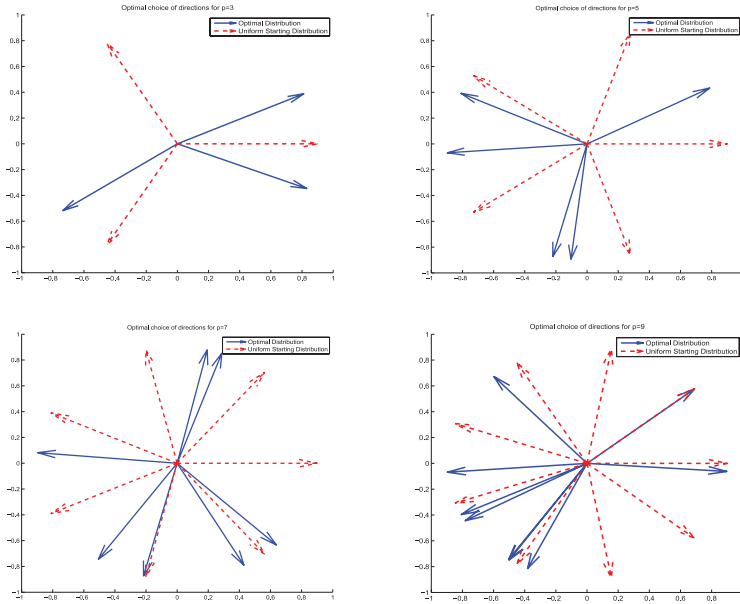


Figure 7. Example 2: Uniform initial distribution θ_0 of the plane wave directions (red dotted line) and the computed optimal distribution (blue solid line) for $p = 3$ (top left), $p = 5$ (top right), $p = 7$ (bottom left), and $p = 9$ (bottom right). Parameter choice (6.1).

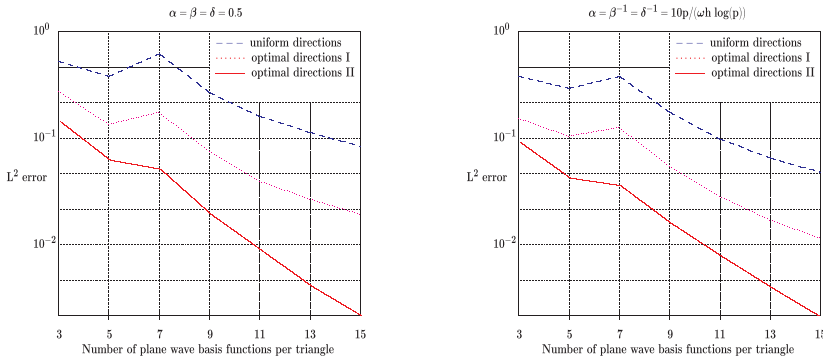


Figure 8. Example 2: The L^2 -error $\|u - u_h\|_{0,\Omega}$ as a function of the number p of plane wave basis functions per element: Uniformly distributed directions (blue broken line), optimal directions I (red dotted line), and optimal directions II (red solid line). Parameter choice (6.1) (left) and parameter choice (6.1) (right).

In case of optimal directions I, Figure 7 shows the uniformly chosen initial distribution θ_0 and the computed optimal distribution of the plane wave directions for $p = 3, 5, 7$, and $p = 9$ (from top left to bottom right). The dependence of the L^2 -norm of the global discretization error $\|u - u_h\|_{0,\Omega}$ is shown in Figure 8 for both optimal directions I and optimal directions II. Since in this example there are no dominant global directions of propagation, the convergence is better for different directions per element than for the same directions.

Acknowledgments. Both authors acknowledge support by the NSF grants DMS-1216857 and DMS-1520886. The second author further acknowledges support by the German National Science Foundation within the Priority Programs SPP 1506 and by the European Science Foundation (ESF) within the ESF Program OPTPDE.

References

- [1] M. Amara, S. Chaudhry, J. Diaz, R. Djellouli, and S. L. Fiedler, A local wave tracking strategy for efficiently solving mid- and high-frequency Helmholtz problems. *Comput. Methods Appl. Mech. Engrg.* **276**, 473–508, 2014.
- [2] M. Amara, R. Djellouli, and C. Farhat, Convergence analysis of a discontinuous Galerkin method with plane waves and lagrange multipliers for the solution of Helmholtz problems. *SIAM J. Numer. Anal.* **47**, 1038–1066, 2009.
- [3] G. B. Alvarez, A. F. D. Loula, E. G. Dutra do Carmo, and F. A. Rochinha, A discontinuous finite element formulation for Helmholtz equation. *Comput. Methods Appl. Mech. Engrg.* **195**, 4018–4035, 2006.
- [4] O. Cessenat and B. Després, Application of an ultra weak variational formulation of elliptic PDEs to the two-dimensional Helmholtz equation. *SIAM J. Numer. Anal.* **35**, 255–299, 1998.
- [5] E. T. Chung and B. Engquist, Optimal discontinuous Galerkin methods for wave propagation. *SIAM J. Numer. Anal.* **44**, 2131–2158, 2006.
- [6] X. Feng and H. Wu, Discontinuous Galerkin methods for the Helmholtz equation with large wave numbers. *SIAM J. Numer. Anal.* **47**, 2872–2896, 2009.
- [7] X. Feng and H. Wu, hp-discontinuous Galerkin methods for the Helmholtz equation with large wave number. *Math. Comp.* **80**, 1997–2024, 2011.
- [8] G. Gabard, Discontinuous Galerkin methods with plane waves for time-harmonic problems. *J. Comp. Phys.* **225**, 1961–1984, 2007.
- [9] C. Gittelsohn, R. Hiptmair, and I. Perugia, Plane wave discontinuous Galerkin methods: Analysis of the h-version. *ESAIM: M2AN Math. Model. Numer. Anal.* **43**, 297–331, 2009.
- [10] R. Griesmaier and P. Monk, Error analysis for a hybridizable discontinuous Galerkin method for the Helmholtz equation. *J. Sci. Comp.* **49**, 291–310, 2011.

- [11] R. Hiptmair, A. Moiola and I. Perugia, Plane wave discontinuous Galerkin methods for the 2D Helmholtz equation: Analysis of the p-version. *SIAM J. Numer. Anal.* **49**, 264–284, 2011.
- [12] R. Hiptmair, A. Moiola and I. Perugia, Plane wave discontinuous Galerkin methods: Exponential convergence of the hp-version. *Found. Comp. Math.* **16**, 637–675, 2016.
- [13] R. H. W. Hoppe and N. Sharma, Convergence analysis of an adaptive Interior Penalty Discontinuous Galerkin method for the Helmholtz equation. *IMA J. Numer. Anal.* **33**, 747–770, 2013.
- [14] K. Ito and K. Kunisch, Lagrange Multiplier Approach to Variational Problems and Applications. SIAM, Philadelphia, 2008.
- [15] P. Monk and D. Q. Wang, A least-squares method for the Helmholtz equation. *Comput. Meths. Appl. Mech. Engrg.* **175**, 121–136, 1999.
- [16] J. Nocedal and S. Wright, Numerical Optimization. 2nd Edition. Springer, Berlin-Heidelberg-New York, 2006.
- [17] Y. Saad and M. H. Schultz, GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM J. Sci. Statist. Comput.* **7**, 856–869, 1986.
- [18] L. Tartar, Introduction to Sobolev Spaces and Interpolation Theory. Springer, Berlin-Heidelberg-New York, 2007.
- [19] F. Tröltzsch, Optimal Control of Partial Differential Equations: Theory, Methods and Applications. American Mathematical Society, Providence, 2010.

Received September 21, 2016; revision received January 24, 2017

A. Agrawal, Department of Bioinformatics, University of Texas at El Paso, El Paso, TX 79968-0766, USA

E-mail: aagrawal@utep.edu

R. H. W. Hoppe, Department of Mathematics, University of Houston, Houston, TX 77204-3008, USA; Institute of Mathematics, University of Augsburg, D-86159 Augsburg, Germany

E-mail: rohop@math.uh.edu; hoppe@math.uni-augsburg.de