# Optimal solvers for PDE-constrained optimization

Tyrone Rees H. Sue Dollar Andrew J. Wathen

2022.11.18

# Contents

## the distributed control problem

we consider the distributed control problem which consists of a cost functional (1) to be minimized subject to a partial differential equation problem posed on a domain $\Omega$ $\mathbb{R}^2$or $\mathbb{R}^3$:

$$\min_{u,f} \frac{1}{2}||u - \hat{u}||_2^2 + \beta||f||_2^2 \qquad (1)$$

$$\text{subject to } -\nabla^2 u = f \quad \text{in } \Omega. \qquad (2)$$

$$\text{with } u = g \quad \text{on } \partial\Omega_1 \quad \text{and} \quad \frac{\partial u}{\partial n} = g \quad \text{on } \partial\Omega_2. \qquad (3)$$

where $\partial\Omega_1 \cup \partial\Omega_2 = \partial\Omega$ and $\partial\Omega_1$ and $\partial\Omega_2$ are distinct. Here, the function $\hat{u}$ (the 'desired state') is known, and we want to find u which satisfies the PDE problem and is as close to $\hat{u}$ as possible in the $L_2$ norm sense.

## Formulation and Structure

For definiteness and clarity we describe this for the purely Dirichlet problem; the formulation for the mixed and purely Neumann problem is also standard. The Dirichlet problem is: find $u \in H_g^1(\Omega) = \{u : u \in H^1(\Omega), u = g \text{ on } \partial\Omega\}$ such that

$$\int_\Omega \nabla u \cdot \nabla v = \int_\Omega vf \qquad \forall v \in H_0^1(\Omega). \tag{4}$$

We assume that $V_0^h \subset H_0^1$ is an n-dimensional vector space of test functions with $\{\phi_1, \cdots, \phi_n\}$ as a basis. Then, for the boundary condition to be satisfied, we extend the basis by defining functions $\phi_{n+1}, \cdots, \phi_{n+\partial_n}$ and coefficients $U_j$ so that $\sum_{j=n+1}^{n+\partial_n} U_j\phi_j$ interpolates the boundary data. Then, if $u_h \in V_g^h \subset H_g^1(\Omega)$, it is uniquely determined by $u = (U_1 \cdots U_n)^T$ in

$$u_h = \sum_{j=1}^n U_j\phi_j + \sum_{j=n+1}^{n+\partial_n} U_j\phi_j.$$

Here the $\phi_i$, $i = 1, \cdots, n$, define a set of shape functions.

## Formulation and Structure

We also assume that this approximation is conforming, i.e.
$V_g^h = \mathrm{span}\{\phi_1, \cdots, \phi_{n+\partial_n}\} \subset H_g^1(\Omega)$. Then we get the finite-dimensional analogue of (4): find $u_h \in V_g^h$ such that

$$\int_\Omega \nabla u_h \cdot \nabla v_h = \int_\Omega v_h f \qquad \forall v_h \in V_0^h.$$

We discretize $f$ using the same basis used for $u$, so

$$f_h = \sum_{j=1}^n F_j \phi_j$$

since it is well known that $f_h = 0$ on $\partial\Omega$. Thus we can write the discrete analogue of the minimization problem as

$$\min_{u_h, f_h} \frac{1}{2} ||u_h - \hat{u}||_2^2 + \beta ||f_h||_2^2 \tag{5}$$

$$\text{such that } \int_\Omega \nabla u_h \cdot \nabla v_h = \int_\Omega v_h f_h \qquad \forall v_h \in V_0^h. \tag{6}$$

## Formulation and Structure

We can write the discrete cost functional as

$$\min_{u_h, f_h} \frac{1}{2}||u_h - \hat{u}||_2^2 + \beta ||f_h||_2^2 = \min_{\mathbf{u}, \mathbf{f}} \frac{1}{2}\mathbf{u}^T M \mathbf{u} - \mathbf{u}^T \mathbf{b} + \alpha + \beta \mathbf{f}^T M \mathbf{f} \qquad (7)$$

where $\mathbf{u} = (U_1, \cdots, U_n)^T, \mathbf{f} = (F_1, \cdots, F_n)^T, \mathbf{b} = \{\int \hat{u}\phi_i\}_{i=1,\cdots,n}, \alpha = ||\hat{u}||_2^2$ and $M = \{\int \phi_i\phi_j\}_{i,j=1,\cdots,n}$ is a mass matrix. We now turn our attention to the constraint: (6) is equivalent to finding $u$ such that

$$\int_\Omega \nabla\bigg(\sum_{i=1}^n U_i\phi_i\bigg)\cdot\nabla\phi_j + \int_\Omega \nabla\bigg(\sum_{i=n+1}^{n+\partial_n} U_i\phi_i\bigg)\cdot\nabla\phi_j = \int_\Omega \bigg(\sum_{i=1}^n F_i\phi_i\bigg)\phi_j, \quad j = 1, \cdots, n$$

which is

$$\sum_{i=1}^n U_i \int_\Omega \nabla\phi_i \cdot \nabla\phi_j = \sum_{i=1}^n F_i \int_\Omega \phi_i\phi_j - \sum_{i=n+1}^{n+\partial_n} U_i \int_\Omega \nabla\phi_i \cdot \nabla\phi_j, \quad j = 1, \cdots, n$$

or

$$K\mathbf{u} = M\mathbf{f} + \mathbf{d} \tag{8}$$

where the matrix $K = \{\int \nabla \phi_i \cdot \nabla \phi_j\}_{i,j=1\cdots n}$ is the discrete Laplacian (the stiffness matrix) and $\mathbf{d}$ contains the terms coming from the boundary values of $u_h$. Thus (7) and (8) together are equivalent to (5) and (6). One way to solve this minimization problem is by considering the Lagrangian

$$\mathcal{L} := \frac{1}{2}\mathbf{u}^T M \mathbf{u} - \mathbf{u}^T \mathbf{b} + \alpha + \beta \mathbf{f}^T M \mathbf{f} + \lambda^T (K\mathbf{u} - M\mathbf{f} - \mathbf{d}),$$

where $\lambda$ is a vector of Lagrange multipliers. Using the stationarity conditions of $\mathcal{L}$, we find that $\mathbf{f}, \mathbf{u}$ and $\lambda$ are defined by the linear system

$$\begin{bmatrix} 2\beta M & 0 & -M \\ 0 & M & K^T \\ -M & K & 0 \end{bmatrix} \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \\ \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{b} \\ \mathbf{d} \end{bmatrix}. \tag{9}$$

Note that this system of equations has saddle-point system structure, i.e. it is of the form

$$\begin{bmatrix} A & B^T \\ B & -C \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}, \tag{10}$$

where $A = \begin{bmatrix} 2\beta M & 0 \\ 0 & M \end{bmatrix}$, $B = \begin{bmatrix} -M & K \end{bmatrix}$, $C = 0$.

This system is usually very large—each of the blocks $K$ is itself a discretization of the PDE—and sparse, since as well as the zero blocks, $K$ and $M$ are themselves sparse because of the finite element discretization.

# Contents

In general, the system (9) will be symmetric but indefinite, so we may use the MINRES algorithm to solve the system: this is a Krylov subspace method for symmetric linear systems. We want to find a matrix (or a linear process) $\mathcal{P}$ for which $\mathcal{P}^{-1}\mathcal{A}$ has better spectral properties (and such that $\mathcal{P}^{-1}\mathbf{v}$ is cheap to evaluate for any given vector $\mathbf{v}$). We then solve a symmetric preconditioned system equivalent to

$$\mathcal{P}^{-1}\mathcal{A}\mathbf{x} = \mathcal{P}^{-1}\mathbf{b}$$

The aim of preconditioning is to choose a matrix $\mathcal{P}$ such that the eigenvalues of $\mathcal{P}^{-1}\mathcal{A}$ are clustered.

# Block Diagonally Preconditioned MINRES

## Theorem 2.1

*If*

$$\mathcal{A} = \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix}$$

*is preconditioned by*

$$\mathcal{P} = \begin{bmatrix} A & 0 \\ 0 & BA^{-1}B^T \end{bmatrix}$$

*Then the preconditioned matrix $\mathcal{T} = \mathcal{P}^{-1}\mathcal{A}$ satisfies*

$$\mathcal{T}(\mathcal{T} - I)(\mathcal{T}^2 - \mathcal{T} - I) = 0.$$

This shows us that $\mathcal{T}$ is diagonalizable and has at most four distinct eigenvalues $(0, 1, \frac{1 \pm \sqrt{5}}{2})$ or only the three non-zero eigenvalues if $\mathcal{T}$ is nonsingular. This means that the Krylov subspace $\mathcal{K}(\mathcal{T}; \mathbf{r}) = span(\mathbf{r}, \mathcal{T}\mathbf{r}, \mathcal{T}^2\mathbf{r}, \cdots)$ will be of dimension at most three if $\mathcal{T}$ is nonsingular or four if $\mathcal{T}$ is singular. Therefore, any Krylov subspace method with an optimality property (such as MINRES) will terminate in at most three iterations (with exact arithmetic).

# Block Diagonally Preconditioned MINRES

If we apply this approach to the matrix in our saddle-point system (9) then we obtain the preconditioner

$$\mathcal{P}_{MGW} = \begin{bmatrix} 2\beta M & 0 & 0 \\ 0 & M & 0 \\ 0 & 0 & \frac{1}{2\beta}M + KM^{-1}K^T \end{bmatrix}$$

MINRES with this preconditioner will always terminate (in exact arithmetic) in at most three steps and so satisfies one requirement of a preconditioner.

The difficulty comes from the (3,3) block, which is the only part that contains the PDE. One way to approximate this is to consider only the dominant term in the (3,3) block which is, for all but the very smallest values of $\beta$, the $KM^{-1}K^T$ term, thus forming the preconditioner:

$$\mathcal{P}_{D1} = \begin{bmatrix} 2\beta M & 0 & 0 \\ 0 & M & 0 \\ 0 & 0 & KM^{-1}K^T \end{bmatrix}$$

## Proposition 2.1

Let $\lambda$ be an eigenvalue of $\mathcal{P}_{D1}^{-1}\mathcal{A}$. Then either $\lambda = 1$,
$\frac{1}{2}(1 + \sqrt{1 + 4\sigma_1}) \leq \lambda \leq \frac{1}{2}(1 + \sqrt{1 + 4\sigma_m})$ or
$\frac{1}{2}(1 + \sqrt{1 - 4\sigma_m}) \leq \lambda \leq \frac{1}{2}(1 - \sqrt{1 + 4\sigma_1})$ where $0 \leq \sigma_1 \leq \cdots \leq \sigma_m$ are the
eigenvalues of $\frac{1}{2\beta}(KM^{-1}K^T)^{-1}M + I$.

# Block Diagonally Preconditioned MINRES

In our tests, we have discretized the problem (1) using bilinear quadrilateral **Q**1 finite elements, and for this choice one can prove the following.

### Proposition 2.2

*Let $\lambda$ be an eigenvalue of $\mathcal{P}_{D1}^{-1}\mathcal{A}$. Then $\lambda$ satisfies one of*

$$\lambda = 1,$$

$$\frac{1}{2}(1 + \sqrt{5 + \frac{2\alpha_1 h^4}{\beta}}) \leq \lambda \leq \frac{1}{2}(1 + \sqrt{5 + \frac{2\alpha_2}{\beta}})$$

$$or \; \frac{1}{2}(1 - \sqrt{5 + \frac{2\alpha_1 h^4}{\beta}}) \leq \lambda \leq \frac{1}{2}(1 - \sqrt{5 + \frac{2\alpha_2}{\beta}})$$

*where $\alpha_1, \alpha_2$ are positive constants independent of h.*

# Contents

# Constraint Preconditioning

The preconditioner takes the form

$$\mathcal{P} = \begin{bmatrix} G & B^T \\ B & 0 \end{bmatrix}$$

where $G \in \mathbb{R}^{l \times l}$ is a symmetric matrix. Let $Z \in \mathbb{R}^{l \times (l-k)}$ be such that its columns span the nullspace of $B$. The PPCG method can be reliably used if both $Z^T A Z$ and $Z^T G Z$ are positive definite. The basic principles behind the PPCG method are as follows. Let $W \in \mathbb{R}^{l \times k}$ be such that the columns of $W$ together with the columns of $Z$ span $\mathbb{R}^l$ and any solution $x^*$ in (10) can be written as

$$x^* = W x_w^* + Z x_z^*. \tag{11}$$

Substituting (11) into (10) and premultiplying the resulting system by
$\begin{bmatrix} W^T & 0 \\ Z^T & 0 \\ 0 & I \end{bmatrix}$, we obtain the linear system

$$\begin{bmatrix} W^T A W & W A Z & W B^T \\ Z^T A W & Z^T A Z & 0 \\ B W & 0 & 0 \end{bmatrix} \begin{bmatrix} x_w^* \\ x_z^* \\ y \end{bmatrix} = \begin{bmatrix} W^T \mathbf{c} \\ Z^T \mathbf{c} \\ \mathbf{d} \end{bmatrix}.$$

Therefore, we may compute $x_w^*$ by solving

$$BWx_w^* = \mathbf{d},$$

and, having found $x_*^w$, we can compute $x_*^z$ by applying the PCG method to the system

$$A_z z x_*^z = c_z$$

where

$$A_z z = Z^T A Z$$

$$c_z = Z^T(\mathbf{c} - A W x_w^*).$$

# Constraint Preconditioning

The following theorem gives the main properties of the preconditioned matrix $\mathcal{P}^{-1}\mathcal{A}$ : the proof can be found in [1].

## Theorem 3.1

Let $\mathcal{A} = \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix}$ and $\mathcal{P} = \begin{bmatrix} G & B^T \\ B & 0 \end{bmatrix}$, where $B \in \mathbb{R}^{k \times l}$ has full rank, $G \in \mathbb{R}^{l \times l}$ is symmetric and $\mathcal{P}$ is nonsingular. Let the columns of $Z \in \mathbb{R}^{l \times (l-k)}$ span the nullspace of B, then $\mathcal{P}^{-1}\mathcal{A}$ has

- *2k eigenvalues at 1; and*
- *the remaining eigenvalues satisfy the generalized eigenvalue problem*
$$Z^T A Z x_z = \lambda Z^T G Z x_z \tag{12}$$

# Constraint Preconditioning

Additionally, if $G$ is nonsingular, then the eigenvalues defined by (12) interlace the eigenvalues of $G^{-1}A$. Keller *etal.* also show that the Krylov subspace

$$\mathcal{K}(\mathcal{P}^{-1}\mathcal{A}; \mathbf{r}) = \mathrm{span}(\mathbf{r}, \mathcal{P}^{-1}\mathcal{A}\mathbf{r}, (\mathcal{P}^{-1}\mathcal{A})^2\mathbf{r}, \cdots)$$

will be of dimension at most $l - k + 2$, see [1].

Clearly, for our problem (9), $A$ is positive definite and, hence, $Z^T A Z$ is positive definite. It remains for us to show that we can choose a symmetric matrix $G$ that satisfies the following properties:

- $Z^T G Z$ is positive definite;
- the eigenvalues of $\mathcal{P}^{-1}\mathcal{A}$ are clustered; and
- we can efficiently carry out solves with $\mathcal{P}$.

# Constraint Preconditioning

- $G = \operatorname{diag}(A)$
- $\mathcal{P}_{C1} = \begin{bmatrix} 0 & 0 & -M \\ 0 & 2\beta K^T M^{-1} K & K^T \\ -M & K & 0 \end{bmatrix}$

Carsten Keller, Nicholas I. M. Gould, and Andrew J. Wathen. Constraint preconditioning for indefinite linear systems. SIAM Journal on Matrix Analysis and Applications, 21(4):1300–1317, 2000.

**Thanks for your attention!**