

Operator preconditioning for a class of inequality constrained optimal control problems

Anton Schiela* & Stefan Ulbrich †

January 9, 2014

Abstract

We propose and analyze two strategies for preconditioning linear operator equations that arise in PDE constrained optimal control in the framework of conjugate gradient methods. Our particular focus is on control or state constrained problems, where we consider the question of robustness with respect to critical parameters. We construct a preconditioner that yields favorable robustness properties with respect to critical parameters.

AMS MSC 2000: 65F08, 49M05, 65J10

Keywords: preconditioner, optimal control, control constraints, state constraints

1 Introduction

In this paper we are concerned with the solution of optimization problems, subject to partial differential equations and inequality constraints on the control and/or the state. Such problems can be considered as optimization problems in infinite dimensional function spaces, and in recent years, algorithms have been constructed which tackle these problems in function space. The common feature of these algorithms is that they can be formulated and analyzed in the infinite dimensional setting, and each step of such an algorithm requires the solution of an infinite dimensional problem. Taking, for example, Newton methods, this means that in each iteration a linear operator equation is solved. In general terms a perturbed saddle point problem of the form

$$\begin{pmatrix} M^*M & A^* \\ A & -CC^* \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix} \quad (1)$$

has to be solved in each Newton step (we will give a derivation and a precise functional analytic setting in the next sections). Of course, implementations have to deal with discretized versions of these subproblems, but have the conceptual advantage that the methods inherit much of the structure of the infinite dimensional problem.

In this paper we pursue this line of thought one step further and construct a preconditioned iterative solver for the linear systems that occur in certain instances of PDE constrained optimization, including control constrained problems and regularizations of state

*Supported by the DFG Research Center MATHEON "Mathematics for key technologies"

†Supported by the DFG within the Graduate School "Computational Engineering", the SPP 1253 "Optimization with PDE Constraints" and by SFB 666 and SFB 805

constrained problems. In particular we consider two block preconditioners for the cg-method in function space, applied to this problem class. One of them is straightforward and already well known. The second is applicable under certain circumstances, namely that the control and the observation take place in the same space, and yields increased robustness with respect to certain critical parameters, which may become small or large, if (1) comes from constrained optimal control problems. It can be viewed as a generalization of a preconditioner proposed in [28] for problems without inequality constraints, which yields bounded condition numbers independent of critical parameters. While complete independence is not possible in the presence of inequality constraints, it turns out that the performance of our preconditioner depends on critical parameters in a very moderate way.

Finally, we want to point out that our results are valid also in discretized settings, where the usually infinite dimensional spaces are replaced by finite dimensional (finite element) subspaces. Our analysis includes but does not require infinite dimensional spaces.

Preconditioning and multigrid for optimality systems in PDE constrained optimization is an active topic of research and there are several lines of research. Early attempts were made by Battermann et al [1, 2]. Borzi [3, 4] considers collective Gauss-Seidel smoothers, while Zulehner et al [23, 28] and Wathen et al [18, 7] propose and analyze block preconditioners for such systems all-at-once multigrid preconditioners are considered in [22]. While cases without inequality constraints are well understood meanwhile, the case of control and/or state constraints is still mostly open. First approaches were taken by Herzog and Sachs [10] which observed lack of robustness of standard block preconditioners in particular for state constrained optimization problems and in [24]. In a very recent preprint [17] a preconditioner with favorable stability properties for state constrained problems was proposed, but the analysis does not provide useful estimates for the condition number. This preconditioner fits into our general framework, which works for control constraints and for state constraints.

2 Theoretical framework

In this section we introduce a set of general assumptions imposed on the system (1), discuss its solution in function space by a preconditioned conjugate gradient or MINRES method, and present a couple of examples that fall into our framework.

2.1 A class of optimal control problems and corresponding saddle point problems

Let us consider as an example the optimal control problem

$$\min \frac{1}{2} \|My - y_d\|_H^2 + \frac{\alpha}{2} \|u\|_U^2 \text{ s.t. } Ay - Bu = 0 \quad u \geq 0. \quad (2)$$

It can be shown that the minimizer (y_*, u_*) of our optimal control problem (2) can be characterized by the following control reduced optimality system, where p_* is the so called *adjoint state* and the optimal control is given pointwise by $u_* = \max\{\alpha^{-1}B^*p_*, 0\}$:

$$\begin{aligned} 0 &= M^*(My_* - y_d) + A^*p_* \\ 0 &= Ay_* - B \max\{\alpha^{-1}B^*p_*, 0\}. \end{aligned} \quad (3)$$

This problem can be solved for example by a semi-smooth Newton method with iteration

variables (y_k, p_k) [11, 26], whose Jacobian matrix can be written as

$$\begin{pmatrix} M^*M & A^* \\ A & -B\alpha^{-1}\chi_{\mathcal{I}}(p_k)B^* \end{pmatrix}, \quad (4)$$

where $\chi_{\mathcal{I}}(p_k) = 0$ for $p_k > 0$ and $\chi_{\mathcal{I}}(p_k) = 1$ for $p_k \leq 0$ in a pointwise sense. In the case of bilateral constraints $\underline{u} \leq u \leq \bar{u}$ one obtains similar systems with $\max\{\alpha^{-1}p_k, 0\}$ replaced by $\text{Proj}_{[\underline{u}, \bar{u}]}(\alpha^{-1}p_k)$ and $\chi_{\mathcal{I}} = 0$ for $\alpha^{-1}p_k \in]\underline{u}, \bar{u}[$ and $\chi_{\mathcal{I}}(p_k) = 1$ otherwise.

To complete our notational framework, we consider p_k fixed, write $\chi_{\mathcal{I}} = \chi_{\mathcal{I}}(p_k)$ and set $C := B\alpha^{-1/2}\chi_{\mathcal{I}}$. Thus, we end up with solving systems of equations of the form:

$$\begin{pmatrix} M^*M & A^* \\ A & -CC^* \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix}. \quad (5)$$

Let us fix the theoretical framework for this system. The following abstract and very basic assumptions will be used throughout the paper.

Assumption 2.1. (Basic Assumptions)

- (i) Assume that the state space Y and the space of adjoints P are reflexive Banach spaces. Assume that the control space U and the space of observations H are Hilbert spaces. Further, as a matter of notation, we use $v^*(v)$ for dual pairings, while $\langle v_1, v_2 \rangle_V$ denotes a Hilbert space scalar product.
- (ii) Let $A : Y \rightarrow P^*$, the differential operator for the state equation, be an isomorphism, which implies that its Banach-space adjoint $A^* : P \rightarrow Y^*$ is an isomorphism as well.
- (iii) Let $C : U \rightarrow P^*$ be a continuous operator, and $M : Y \rightarrow H$ a continuous operator with dense range.

We denote by the adjoint $C^* : P \rightarrow U$ the mapping that satisfies

$$\langle C^*p, u \rangle_U = (Cu)(p) \quad \forall u \in U$$

It is continuous as well. Analogously, the adjoint $M^* : H \rightarrow Y^*$, defined via

$$(M^*h)(y) = \langle My, h \rangle_H \quad \forall y \in Y.$$

It is continuous and injective since M has dense range.

Remark 2.2. Alternatively, one could define $C^* : P \rightarrow U^*$ via $(C^*p)(u) = (Cu)(p)$, and analogously M^* . This, however, would necessitate to introduce the Riesz isomorphisms R_U of U and R_M of M into the notation. Since in all our applications, U and M are of the form $L_2(Q)$, these R_U (and analogously R_M) are simply the canonical representation

$$(R_U u)(v) = \int_Q uv \, d\omega \quad \forall v \in U.$$

For this reason we decided to use the definitions of Assumption 2.1(iii) for M^* and C^* leading to a more concise notation. Equivalently, one could introduce the convention that U and H are identified with their duals.

With these definitions our system of equations (5) is just another way of writing down the following weak form:

$$\begin{aligned} \langle M\delta y, Mv \rangle_H + (Av)(\delta p) &= f(v) & \forall v \in Y \\ (A\delta y)(w) - \langle C^*\delta p, C^*w \rangle_U &= g(w) & \forall w \in P. \end{aligned}$$

The reduced system formulation (3) has several advantages, compared to a classical KKT system, containing y, p, u and additional Lagrangian multipliers for the constraint $u \geq 0$. First of all, it is a system of two PDEs, and thus, the solutions of these system are contained in a smoother space than the corresponding right hand sides. This has fundamental consequences for the convergence theory of Newton's method, applied to this system [19]. Second, the use of a projection formula makes an additional special treatment (such as barrier of penalty regularization) of the control constraints unnecessary, and the system can be solved directly. Third, only the smooth variables y and p have to be discretized, leading to optimal discretization schemes [13].

If we apply Galerkin's method to discretize (4), then the same definitions as above can be made on finite dimensional subspaces $Y_h \subset Y$ and $P_h \subset P$. Moreover, in the following derivations, no mesh-size parameter appears, so that the following results are automatically independent of the choice of the mesh.

2.2 Reduction to a convex minimization problem

Next we are going to define a bilinear form $\langle \cdot, \cdot \rangle_K$ together with its domain of definition, which we will denote by $D_K \subset P$. In general D_K is a proper subset of P . This requires a slight shift in our functional analytic framework from the spaces used in Assumption 2.1, which form a natural setting for PDEs in variational form, to modified spaces, more suitable for the analysis of our saddle point system. This will take place in two steps.

As a first step, we will introduce a stronger norm (in fact a scalar product) on $\text{ran } M^*$. Density of $\text{ran } M$ in H (cf. Assumption 2.1(iii)) implies injectivity of $M^* : H \rightarrow Y^*$, so that $M^* : H \rightarrow \text{ran } M^*$ can be considered as a bijective operator with inverse $M^{-*} := (M^*)^{-1}$. By definition of the scalar product

$$\langle v, w \rangle_{M^*} := \langle M^{-*}v, M^{-*}w \rangle_H \quad \forall v, w \in \text{ran } M^*$$

the induced norm $\|\cdot\|_{M^*}$ renders M^* an isomorphism by construction. By continuity of M^* this norm is stronger than $\|\cdot\|_{Y^*}$, and $(\text{ran } M^*, \langle \cdot, \cdot \rangle_{M^*})$ inherits the Hilbert space structure from H , because as already mentioned $M^* : H \rightarrow (\text{ran } M^*, \langle \cdot, \cdot \rangle_{M^*})$ is an isomorphism.

Let us denote the dual space $(\text{ran } M^*, \|\cdot\|_{M^*})^*$ by \widehat{Y} so that $\widehat{Y}^* = (\text{ran } M^*, \|\cdot\|_{M^*})$ by reflexivity. Since $\|\cdot\|_{M^*}$ is stronger than $\|\cdot\|_{Y^*}$ the embedding $\widehat{Y}^* \hookrightarrow Y^*$, induced by the inclusion $\text{ran } M^* \subset Y^*$ is continuous and injective. Thus its adjoint operator maps $Y = Y^{**}$ continuously into $\widehat{Y} = \widehat{Y}^{**}$ and has dense range.

Remark 2.3. The mapping $Y \rightarrow \widehat{Y}$ is, however, not injective in general. This is only true, if $\text{ran } M^*$ is dense in Y^* , which in turn holds, if $M : Y \rightarrow H$ is injective. If that is the case, \widehat{Y} can be interpreted as an *extension* of Y . Otherwise, \widehat{Y} is rather an extension of $Y/\ker M$ and thus a space of equivalence classes of functions.

As a second step we define the new space

$$D_K := A^{-*}(\text{ran } M^*) = \{p \in P : A^*p \in \text{ran } M^*\} \subset P.$$

Then $M^{-*}A^*p$ is well defined for all $p \in D_K$ and $M^{-*}A^*$ is a bijective mapping $D_K \rightarrow H$. Thus, on D_K the following bilinear form is well defined and positive definite:

$$\langle v, w \rangle_K := \langle M^{-*}A^*v, M^{-*}A^*w \rangle_H + \langle C^*v, C^*w \rangle_U,$$

and D_K is the domain of definition of $\langle \cdot, \cdot \rangle_K$. Since $A^* : P \rightarrow Y^*$ is an isomorphism, and $M^* : H \rightarrow Y^*$ is continuous, the induced norm $\| \cdot \|_K$ is *stronger* than $\| \cdot \|_P$ (but usually not equivalent).

Thus, considering D_K equipped with the scalar product $\langle \cdot, \cdot \rangle_K$ it is easy to see that $M^{-*}A^* : (D_K, \langle \cdot, \cdot \rangle_K) \rightarrow H$ is an isomorphism, taking into account the continuity of $C^* : P \rightarrow U$. Hence, $(D_K, \langle \cdot, \cdot \rangle_K)$ inherits the Hilbert space structure from H .

So, we can consider the following minimization problem for $\ell \in D_K^*$,

$$\psi(p) := \frac{1}{2} \langle p, p \rangle_K + \ell(p), \quad \min_{p \in D_K} \psi(p) \quad (6)$$

Since $(D_K, \langle \cdot, \cdot \rangle_K)$ is a Hilbert space, this problem has a unique solution \tilde{p} and the mapping $\ell \rightarrow \tilde{p}$ is nothing other than the Riesz isomorphism in D_K .

Since $(D_K, \langle \cdot, \cdot \rangle_K)$ is continuously embedded into P , we conclude by duality, that the restriction mapping $P^* \rightarrow D_K^*$ is also continuous, so that each continuous linear functional on P is also continuous on D_K . However, the mapping $P^* \rightarrow D_K^*$ is not injective in general. This is only true if D_K is dense in P , which in turn only holds, if M is injective.

Lemma 2.4. *The minimization problem (6) has a unique solution in D_K for any $\ell \in D_K^*$. Moreover, for any right hand sides $f \in Y^*$ and $g \in P^*$ the system (5) has a solution $(\delta y, \delta p) \in Y \times P$. It can be computed from the solution \tilde{p} of*

$$\min_{p \in D_K} \frac{1}{2} \langle p, p \rangle_K + g(p) + \langle C^*A^{-*}f, C^*p \rangle_U. \quad (7)$$

via

$$\begin{aligned} \delta p &= \tilde{p} + A^{-*}f \text{ in } Y^* \\ \delta y &= A^{-1}(g + CC^*\delta p) \text{ in } P^*. \end{aligned} \quad (8)$$

Proof. As already discussed, a unique minimizer of (6) exists, and its first order optimality conditions read

$$\langle M^{-*}A^*\tilde{p}, M^{-*}A^*w \rangle_H + \langle C^*\tilde{p}, C^*w \rangle_U + \ell(w) = 0 \quad \forall w \in D_K.$$

Setting $\ell(p) := g(p) + \langle C^*A^{-*}f, C^*p \rangle_U$ this yields

$$\langle M^{-*}A^*\tilde{p}, M^{-*}A^*w \rangle_H + \langle C^*\tilde{p}, C^*w \rangle_U + g(w) + \langle C^*A^{-*}f, C^*w \rangle_U = 0 \quad \forall w \in D_K. \quad (9)$$

By definition, δy and δp solve the second row of (5), so it remains to show that they solve the first row. Inserting (8) we conclude via $A\delta y = CC^*\tilde{p} + g + CC^*A^{-*}f$ that

$$\langle M^{-*}(A^*\delta p - f), M^{-*}A^*w \rangle_H + (A\delta y)(w) = 0 \quad \forall w \in D_K. \quad (10)$$

For arbitrary $v \in Y$ define w as the solution of the equation $A^*w = M^*Mv$, or more explicitly

$$(A\eta)(w) = \langle M\eta, Mv \rangle_H \quad \forall \eta \in Y.$$

By definition, $A^*w \in \text{ran } M^*$ and thus $w \in D_K$. Then in particular $(A\delta y)(w) = \langle M\delta y, Mv \rangle_H$, and $M^{-*}A^*w = Mv$, and we conclude from (10)

$$\langle M^{-*}(A^*\delta p - f), Mv \rangle_H + \langle M\delta y, Mv \rangle_H = 0 \quad \forall v \in Y$$

which yields

$$(A^*\delta p - f)(v) + (M^*M\delta y)(v) = 0 \quad \forall v \in Y$$

and thus, in short, the first row of (5). \square

Hence, we can find the solution of our block system (5) by solving (7) for \tilde{p} , which is equivalent to solving (9), and then computing δp and δy by (8).

Let us finally represent the first derivative of $\psi(p)$ in operator form. Obviously, we have

$$\psi'(p)\delta p = \langle p, \delta p \rangle_K + \ell(\delta p) = \langle M^{-*}A^*p, M^{-*}A^*\delta p \rangle_H + \langle C^*p, C^*\delta p \rangle_U + \ell(\delta p) \quad \forall \delta p \in D_K.$$

We observe that $M^* : H \rightarrow \widehat{Y}^*$ and $A^* : D_K \rightarrow \widehat{Y}^*$ are isomorphisms by definition of the involved norms, and $C^* : D_K \rightarrow U$ is continuous. Let us denote by $\widehat{M} : \widehat{Y} \rightarrow H$, $\widehat{A} : \widehat{Y} \rightarrow D_K^*$, their respective adjoint operators in the framework of these spaces. Then we can write:

$$\psi'(p)\delta p = (\widehat{A}\widehat{M}^{-1}M^{-*}A^*p + CC^*p + \ell)(\delta p) \quad \forall \delta p \in D_K.$$

The notational distinction has to be done for \widehat{A} and \widehat{M} , since their domain space differs from the one of A and M . This is not necessary for C , so $\widehat{C} = C^{**} = C$ needs not to be distinguished from C .

Remark 2.5. Recall, that if M is injective, the mapping $Y \rightarrow \widehat{Y}$ is injective and dense. Then \widehat{A} and \widehat{M} are extensions of A and M , respectively. Otherwise, one has to factor out $\ker M$, as described in Remark 2.3. Then \widehat{A} and \widehat{M} are mappings in spaces of equivalence classes.

2.3 Preconditioned conjugate gradient method in function space

In this paper we concentrate for simplicity on the solution of (1) via the reduction to the convex problem (6) and by applying a preconditioned conjugate gradient method. Several alternatives have been proposed in the literature, such as preconditioned MINRES [18] and a Bramble-Pasciak cg-method [10]. The application of our preconditioners in the context of the preconditioned MINRES method will be considered in 2.4.

On the Hilbert space $(D_K, \langle \cdot, \cdot \rangle_K)$ consider once more the convex, quadratic functional $\psi(p) = \frac{1}{2}\langle p, p \rangle_K + \ell(p)$ from (6). Further, let a different scalar product $\langle \cdot, \cdot \rangle_Q$ (a preconditioner) be given on D_K . Denote by $\nabla_Q\psi(p)$ the Q -gradient of ψ at p , the unique vector $\nabla_Q\psi(p) \in D_K$ that satisfies

$$\langle \nabla_Q\psi(p), w \rangle_Q = \psi'(p)(w) = \ell(w) + \langle p, w \rangle_K \quad \forall w \in D_K.$$

Then the method of conjugate gradients can be written as follows (cf. also [9]):

Algorithm 2.6. (preconditioned cg in function space)

$$p_0 \text{ given, } d_0 := -\nabla_Q\psi(p_0)$$

$k = 0, 1, 2, \dots$

$$\begin{aligned} p_{k+1} &= p_k - \frac{\psi'(p_k)d_k}{\langle d_k, d_k \rangle_K} d_k \quad (\text{exact linesearch along } d_k) \\ g_{k+1} &= -\nabla_Q \psi(p_{k+1}) \quad (\text{direction of steepest descent w.r.t } \langle \cdot, \cdot \rangle_Q) \\ d_{k+1} &= g_{k+1} - \frac{\langle g_{k+1}, d_k \rangle_K}{\langle d_k, d_k \rangle_K} d_k \quad (\text{orthogonalization w.r.t. } \langle \cdot, \cdot \rangle_K) \end{aligned}$$

In the remainder of the paper we consider application of Algorithm 2.6 to (6) and construct and analyze bilinear forms $\langle \cdot, \cdot \rangle_Q$ on D_K . We will establish estimates of the form (11) such that κ_Q is small. In particular we want to avoid that κ_Q depends strongly on certain critical parameters, e.g. α in Example 2.5.3 or γ in Example 2.5.4, that arise in optimal control problems.

It is well known that speed of convergence of the cg-method depends on the condition number κ_Q of $\langle \cdot, \cdot \rangle_K$ with respect to $\langle \cdot, \cdot \rangle_Q$. It can be defined as follows. If the following (sharp) estimates hold,

$$m_Q \langle v, v \rangle_Q \leq \langle v, v \rangle_K \leq M_Q \langle v, v \rangle_Q \quad \forall v \in D_K, \quad (11)$$

then the condition number is given by $\kappa_Q := M_Q/m_Q$. Then, if \tilde{p} denotes the minimizer of ψ , we have the well known estimate

$$\|p_k - \tilde{p}\|_K \leq 2 \left(\frac{\sqrt{\kappa_Q} - 1}{\sqrt{\kappa_Q} + 1} \right)^k \|p_0 - \tilde{p}\|_K \quad (12)$$

and the number of iterations to reach a certain accuracy is proportional to $\sqrt{\kappa_Q}$ (cf. e.g. [6, Sec. 5.3.2]). Thus, it is crucial to find a good preconditioner Q that renders κ_Q small.

Remark 2.7. In finite dimensions (11) is often formulated as a condition on generalized eigenvalues. If Q and K are the matrices, corresponding to $\langle \cdot, \cdot \rangle_Q$ and $\langle \cdot, \cdot \rangle_K$ respectively, and $\lambda_{\min}, \lambda_{\max}$ are the extreme eigenvalues of the problem $K - \lambda Q = 0$, then

$$\kappa_Q = \frac{\lambda_{\max}}{\lambda_{\min}}.$$

2.4 An alternative: preconditioned MINRES method

While we concentrate on the construction of preconditioners for the reduced system (9), we briefly show how they can be converted to preconditioners for MINRES applied to the full system (5). Instead of solving (9) by a preconditioned conjugate gradient method we can also solve the block system

$$\begin{aligned} & \begin{pmatrix} M^* \widehat{M} & A^* \\ \widehat{A} & -CC^* \end{pmatrix} \begin{pmatrix} \tilde{y} \\ \tilde{p} \end{pmatrix} = \begin{pmatrix} 0 \\ g + CC^* A^{-*} f \end{pmatrix}, \quad (13) \\ B & := \begin{pmatrix} M^* \widehat{M} & A^* \\ \widehat{A} & -CC^* \end{pmatrix} : \widehat{Y} \times D_K \rightarrow \widehat{Y}^* \times D_K^* \end{aligned}$$

with the Hilbert spaces D_K and \widehat{Y} and the isomorphisms $M^* : H \rightarrow \widehat{Y}^*$, $A^* : D_K \rightarrow \widehat{Y}^*$, $\widehat{M} : \widehat{Y} \rightarrow H$ and $\widehat{A} : \widehat{Y} \rightarrow D_K^*$ introduced in 2.2. Then the reduced system (9) can be written as

$$K\tilde{p} = -g - CC^* A^{-*} f$$

with the Schur complement $K = \widehat{A}(M^*\widehat{M})^{-1}A^* + CC^* : D_K \rightarrow D_K^*$. Moreover, we have the factorization

$$B = \begin{pmatrix} M^*\widehat{M} & A^* \\ \widehat{A} & -CC^* \end{pmatrix} = \begin{pmatrix} I & 0 \\ \widehat{A}(M^*\widehat{M})^{-1} & I \end{pmatrix} \begin{pmatrix} M^*\widehat{M} & 0 \\ 0 & -K \end{pmatrix} \begin{pmatrix} I & (M^*\widehat{M})^{-1}A^* \\ 0 & I \end{pmatrix}.$$

Let $Q = S^*S : D_K \rightarrow D_K^*$ with $S : D_K \rightarrow H$ be one of the symmetric positive definite preconditioners for K proposed in this paper. MINRES requires a positive definite preconditioner. Based on the above factorization we consider the choice

$$T = \begin{pmatrix} I & 0 \\ \widehat{A}(M^*\widehat{M})^{-1} & I \end{pmatrix} \begin{pmatrix} M^*\widehat{M} & 0 \\ 0 & Q \end{pmatrix} \begin{pmatrix} I & (M^*\widehat{M})^{-1}A^* \\ 0 & I \end{pmatrix}$$

Then $T : \widehat{Y} \times D_K \rightarrow \widehat{Y}^* \times D_K^*$ is symmetric positive definite and

$$T^{-1} = R^*R, \quad R = \begin{pmatrix} M^{-*} & 0 \\ 0 & S^{-*} \end{pmatrix} \begin{pmatrix} I & 0 \\ -\widehat{A}(M^*\widehat{M})^{-1} & I \end{pmatrix}, \quad (14)$$

where $R : \widehat{Y}^* \times D_K^* \rightarrow H \times H$. By using the above factorization it is easy to see that

$$RBR^* = \begin{pmatrix} I & 0 \\ 0 & -S^{-*}KS^{-1} \end{pmatrix}. \quad (15)$$

Hence, the preconditioned full system has the eigenvalue 1 and the eigenvalues of $-Q^{-1}K$.

Now (13) has the form

$$Bx = b$$

with a symmetric operator $B \in \mathcal{L}(X, X^*)$, i.e. $\langle x, By \rangle = \langle y, Bx \rangle$. Let $T^{-1} = R^*R : X^* \rightarrow X$ be a symmetric and coercive preconditioner with $R \in \mathcal{L}(X^*, H)$ and a Hilbert space H .

The preconditioned MINRES algorithm computes $x_k = R^*y_k$ with

$$\min \|RBR^*y - Rb\|_H \text{ s.t. } y_k \in y_0 + K_k(RBR^*, Rr_0),$$

where $r_0 = Bx_0 - b$ and the Krylov space $K_k(RBR^*, Rr_0) = \text{span}(Rr_0, \dots, (RBR^*)^{k-1}Rr_0)$. We now estimate the speed of convergence. We have

$$\|Rr_k\|_H = \min_{q \in \pi_k} \|q(RAR^*)Rr_0\|_H,$$

where π_k is the set of polynomials of degree $\leq k$ with $q(0) = 1$. Now let $V = (v_1, \dots, v_N)$ be an arbitrary orthonormal system in H with $K_{k+1}(RBR^*, Rr_0) \subset \text{span}(V)$. Then the matrix $C = \langle V, RBR^*V \rangle_H$ is symmetric and we find a matrix Q with $\text{diag}(\lambda_1, \dots, \lambda_N) = Q^T M Q$, $Q^T Q = I$. Hence $W = (w_1, \dots, w_N) = VQ$ satisfies $\langle w_i, w_j \rangle_H = \delta_{ij}$ and $\langle w_i, RBR^*w_j \rangle_H = \delta_{ij}\lambda_i$. Let $Rr_0 = \sum_{j=1}^N z_j w_j$. We know that $q(RBR^*)Rr_0 \in \text{span}(W)$ and since $\langle w_i, RBR^*Rr_0 \rangle_H = \lambda_i z_i$, we conclude that $RBR^*Rr_0 = \sum_{j=1}^N \lambda_j z_j w_j$ and by induction $q(RBR^*)Rr_0 = \sum_{j=1}^N q(\lambda_j) z_j w_j$. Hence, $\|q(RBR^*)Rr_0\|_H \leq \|Rr_0\|_H \max_j |q(\lambda_j)|$ and we have shown that

$$\|Rr_k\|_H \leq \min_{q \in \pi_k} \max_j |q(\lambda_j)| \|Rr_0\|_H,$$

where λ_j are the eigenvalues of the matrix $C = \langle V, RAR^*V \rangle_H$.

If we apply the preconditioned MINRES method with preconditioner (14) to the system (13) we obtain the following result.

Lemma 2.8. *Let $J \subset]0, \infty[$ contain the spectrum of $Q^{-1}K$, more precisely $\frac{\langle v, v \rangle_K}{\langle v, v \rangle_Q} \in J$ for all $0 \neq v \in D_K$. If the preconditioned MINRES method is applied to (13) with preconditioner (14) then it computes $x_k = R^*y_k$ with*

$$\|Rr_k\|_H \leq \min_{q \in \pi_k} \max_{\lambda \in J \cup \{-1\}} |q(\lambda)| \|Rr_0\|_H \leq \min_{q \in \pi_{k-1}} \max_{\lambda \in J} |(\lambda + 1)q(\lambda)| \|Rr_0\|_H, \quad (16)$$

where $r_k = Bx_k - b$.

Proof. RBR^* has the block diagonal form (15). Let $V = \left(\begin{pmatrix} V_1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ V_2 \end{pmatrix} \right)$ be an orthonormal system in $H \times H$ such that $K_{k+1}(RBR^*, Rr_0) \subset \text{span}(V)$. Then the matrix $C = \langle V, RBR^*V \rangle_{H \times H}$ has the form

$$C = \langle V, RBR^*V \rangle_{H \times H} = \begin{pmatrix} I & 0 \\ 0 & -\langle V_2, S^{-*}KS^{-1}V_2 \rangle_H \end{pmatrix}.$$

The eigenvalues of the symmetric matrix $C_2 = \langle V_2, S^{-*}KS^{-1}V_2 \rangle_H$ are contained in J , since for alle $0 \neq z \in \mathbb{R}^{k+1}$ with $w = S^{-1}V_2z$ holds

$$\frac{z^T C_2 z}{z^T z} = \frac{\langle V_2 z, S^{-*}KS^{-1}V_2 z \rangle_H}{\langle V_2 z, V_2 z \rangle_H} = \frac{\langle w, w \rangle_K}{\langle w, w \rangle_Q} \in J.$$

Hence, the estimate (16) is shown, where the second inequality follows from the fact that for all $q \in \pi_{k-1}$ the polynomial $(\cdot + 1)q$ is in π_k and vanishes at -1 . \square

2.5 Examples of optimal control problems

To clarify our abstract setting and our choice of spaces, let us consider several examples for optimal control problems. We give a couple of examples for the abstract operators A, B, C, M and the spaces H, U, P, Y involved. For simplicity of presentation we consider linear-quadratic problems here. Nonlinear problems can be solved iteratively via a Newton type algorithm, which requires the solution of a linear system of operator equations in each step.

2.5.1 Elliptic optimal control problems

For elliptic optimal control on a bounded Lipschitz domain $\Omega \in \mathbb{R}^d$ with boundary Γ , let $H_0^1(\Omega) \subset V \subset H^1(\Omega)$ be an appropriate closed subspace of $H^1(\Omega)$ that incorporates boundary conditions. For example, $V = H_0^1(\Omega)$ corresponds to homogeneous Dirichlet boundary conditions, and $V = H^1(\Omega)$ to Neumann or Robin boundary conditions.

Let $Y := P := V$ and define A as follows

$$A : Y = V \rightarrow P^* = V^* \\ y \mapsto Ay : (Ay)(v) := a(y, v) := \int_{\Omega} \langle \sigma(x) \nabla y, \nabla v \rangle_{\mathbb{R}^d} + a_0(x) y v \, dx, \quad (17)$$

where $\sigma(x) : \Omega \rightarrow \mathbb{R}^{d \times d}$ defines a symmetric bounded elliptic bilinear form which is continuous in x , and $a_0 : \Omega \rightarrow \mathbb{R}$ is nonnegative and bounded. In the case $V = H^1(\Omega)$ we additionally assume that a_0 is positive on a subset Ω of non-zero measure. This ensures that $a(\cdot, \cdot)$ is V -elliptic, i.e., there is a constant $c_a > 0$, such that

$$a(v, v) \geq c_a \|v\|_{H^1}^2 \quad \forall v \in V.$$

Having fixed the framework for $A : Y \rightarrow P^*$, let us consider two variants in the choice of H, U, M, C and B . Other combinations are conceivable, as well.

Distributed control and observation. In this case we set $U = H = L_2(\Omega)$ and $Y = P = H_0^1(\Omega)$. Let us define the operator $M = E_S$ as the Sobolev embedding $E_S : H_0^1(\Omega) \hookrightarrow L_2(\Omega)$, while $B : L_2(\Omega) \rightarrow H_0^1(\Omega)^*$ is, in the case of distributed control defined via

$$(Bu)(v) = \int_{\Omega} u \cdot E_S v, dx = (E_S^* u)(v). \quad (18)$$

We observe that, as required $M : Y \rightarrow H$ is continuous and has dense range. Even more, M is injective, which implies that $\text{ran } M^*$ is dense in Y^* .

For illustration, let us discuss the shift of framework of Section 2.2 in our concrete setting. Here M^* is the embedding of $H = L_2(\Omega)$ into $Y^* = H_0^1(\Omega)^*$, so that $\text{ran } M^*$ can be identified with $L_2(\Omega)$ as a subset of $H_0^1(\Omega)^*$, and $\|\cdot\|_{M^*} = \|M^{-*} \cdot\|_H = \|\cdot\|_{L_2}$. Thus, we simply have $\widehat{Y} = L_2(\Omega)$ and $\widehat{M} = id_{L_2(\Omega)}$.

The definition of D_K in our concrete example is consequently

$$D_K = \{p \in P : \exists f \in L_2(\Omega) : (A^* p)(v) = a(v, p) = \int_{\Omega} f v dx \quad \forall v \in H_0^1(\Omega)\}.$$

In particular, $D_K \supset H^2(\Omega) \cap H_0^1(\Omega)$ and if Ω and σ are sufficiently regular, $D_K = H^2(\Omega) \cap H_0^1(\Omega)$. In any case D_K is dense in P so that the mapping $P^* \rightarrow D_K^*$ is injective.

Remark 2.9. Comparison with the method of transposition (cf. e.g. [16]) yields, that $A^* : D_K \rightarrow \widehat{Y}^*$ can be interpreted as the *strong* form of (17) (as a mapping $p \rightarrow \text{div } \sigma \nabla p + a_0 p := f$ from the solution to the data), while its adjoint, which we called $\widehat{A} : \widehat{Y} = L_2(\Omega) \rightarrow D_K^*$, is an extension of (17), the so called *very weak* form.

Boundary control and observation at the boundary. In this case we set $U = H = L_2(\Gamma)$, and in order to allow Neumann boundary control we set $Y = P = H^1(\Omega)$. Let us define the operator $M = \tau$ as the trace operator $\tau : H^1(\Omega) \hookrightarrow L_2(\Gamma)$, while $B : L_2(\Gamma) \rightarrow H^1(\Omega)^*$ is, in the case of boundary control defined via

$$(Bu)(v) = \int_{\Gamma} u \tau v dS.$$

Also here $M : Y \rightarrow H$ is continuous and has dense range, but now M is not injective, so that $\text{ran } M^*$ is not dense in Y^* .

In our concrete setting M^* is the embedding of $H = L_2(\Gamma)$ into $Y^* = H^1(\Omega)^*$, so that $\text{ran } M^*$ can be identified with $L_2(\Gamma)$ as a subset of $H^1(\Omega)^*$, $\widehat{Y} = L_2(\Gamma)$ and $\widehat{M} = id_{L_2(\Gamma)}$. Obviously, the mapping $Y \rightarrow \widehat{Y}$ is not injective in this case.

Here the definition of D_K reads as follows

$$D_K = \{p \in P : \exists g \in L_2(\Gamma) : (A^* p)(v) = a(v, p) = \int_{\Gamma} g \tau v dS \quad \forall v \in H^1(\Omega)\}.$$

An interpretation is that D_K consists of all functions p with $-\text{div } \sigma \nabla p + a_0 p = 0$ (in the distributional sense) and outer (distributional) normal derivative $\partial_{\sigma\nu} p = g \in L_2(\Gamma)$.

Similar to the case of distributed control $A^* : D_K \rightarrow L_2(\Gamma)$ can be interpreted as a mapping $p \rightarrow \partial_{\sigma\nu} p := g$ from solution to data. The evaluation of $\langle p, w \rangle_K$ for $p, w \in D_K$ is thus performed as

$$\langle p, w \rangle_K = \langle \partial_{\sigma\nu} p, \partial_{\sigma\nu} w \rangle_{L_2(\Gamma)} + \alpha^{-1} \langle \chi_{\mathcal{I}} \tau p, \chi_{\mathcal{I}} \tau w \rangle_{L_2(\Gamma)}.$$

2.5.2 Parabolic optimal control with control constraints

Consider now the parabolic optimal control problem

$$\min \frac{1}{2} \|My - y_d\|_{L_2([0,T], L_2(\Omega))}^2 + \frac{\alpha}{2} \|u\|_{L_2([0,T], L_2(\Omega))}^2 \text{ s.t. } Ay - Bu = 0 \quad u \geq 0.$$

This case runs quite similarly to the elliptic case, with the main difference that

$$A : W([0, T]) \rightarrow L_2([0, T], H^1(\Omega))^*$$

is now a parabolic operator on

$$W([0, T]) = \{y \in L_2([0, T], H^1(\Omega)), y_t \in L_2([0, T], H^1(\Omega)^*)\},$$

defined via

$$(Ay)(v) = \int_{[0,T]} \left(y_t(v) + \int_{\Omega} \nabla y \cdot \sigma(t, x) \nabla v + a_0(t, x) y v \, dx \right) dt.$$

Here $y_t(v)(t)$ is the application of the weak derivative $y_t(t) \in H^1(\Omega)^*$ to $v(t) \in H^1(\Omega)$, which yields an integrable function in time. For a more detailed description consider, e.g., [25]. We use the spaces $Y = W([0, T])$, $P = L_2([0, T], H^1)$, $U = H = L_2([0, T], L_2(\Omega))$ and can now set $M = E_W : Y = W([0, T]) \hookrightarrow H = L_2([0, T], L_2(\Omega))$, the Sobolev embedding for these spaces, and $B : U = L_2([0, T], L_2(\Omega)) \rightarrow P^* = L_2([0, T], H^1(\Omega)^*)$ via

$$(Bu)(v) = \int_{[0,T]} \int_{\Omega} u E_W v \, dx \, dt.$$

Similar to the elliptic case a control reduced optimality system and its semi-smooth Newton linearization can be derived, which is again of the form (4).

In all of the above examples the parameter $\alpha > 0$ appears as a finite (possibly small) but fixed value. The next two examples describe situations, where during algorithmic progress towards the solution such a parameter is adjusted, and makes the problem at hand more and more difficult as the solution is approached.

2.5.3 Regularized bang-bang control

In some application so called bang-bang control is of interest, here written down for distributed control and observation:

$$\min \frac{1}{2} \|My - y_d\|_{L_2(\Omega)}^2 \text{ s.t. } Ay - Bu = 0 \quad \underline{u} \leq u \leq \bar{u}.$$

Usually an optimal solution of such a problem almost everywhere takes either the value \underline{u} or \bar{u} . A simple idea to solve this problem is to consider its regularized versions

$$\min \frac{1}{2} \|My - y_d\|_{L_2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L_2(\Omega)}^2 \text{ s.t. } Ay - Bu = 0 \quad \underline{u} \leq u \leq \bar{u}.$$

and pass to the limit $\alpha \rightarrow 0$ (cf. e.g. [27]). Under certain assumptions on the adjoint state p it can be shown that the regularized solutions u_α tend to the solution of the original problem, and that the Lebesgue measure of the set $\mathcal{I} := \{x \in \Omega : \underline{u} < u_\alpha(x) < \bar{u}\}$ becomes smaller and smaller. In most cases it is observed that $\text{meas}(\mathcal{I}) \leq C\alpha$.

In this setting it is desirable to obtain a preconditioner that is robust for $\alpha \rightarrow 0$.

2.5.4 Elliptic optimal control with state constraints

In the case of state constraints, algorithms typically apply some kind of regularization technique [12, 20], usually based on classical approaches such as penalty or barrier methods. As an example, consider a classical penalty approach, also called “generalized Moreau-Yosida” regularization in this context, for the problem

$$\min \frac{1}{2} \|E_S y - y_d\|_{L_2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L_2(\Omega)}^2 \text{ s.t. } Ay - Bu = 0 \quad y \geq 0,$$

where $E_S : H^1(\Omega) \hookrightarrow L_2(\Omega)$ is again the Sobolev embedding and B , which shall correspond to distributed control, is defined as in (18). Penalization of the constraint $y \geq 0$ yields the problem

$$\min \frac{1}{2} \|E_S y - y_d\|_{L_2(\Omega)}^2 + \frac{\gamma}{2} \|\min\{y, 0\}\|_{L_2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L_2(\Omega)}^2 \text{ s.t. } Ay - Bu = 0,$$

which is often tackled by semi-smooth Newton methods. To approach the solution of the original problem, γ is driven towards $+\infty$ in a path-following method. In practice, the algorithm is terminated with values of γ in the range of 10^8 to 10^{12} . The presence of large γ affects the condition number of the problem severely, if no appropriate measures are taken.

Similarly, this problem can be tackled by barrier methods with a barrier functional $l(\cdot; \mu) :]0, \infty[\rightarrow \mathbb{R}$ parametrized by $\mu > 0$ such that $\lim_{t \rightarrow 0} l(t; \mu) = +\infty$. Path-following algorithms drive μ towards 0 to converge towards the original solution, and similar effects for the condition number occur.

In both cases computing a Newton step amounts in the solution of a linear system of the form

$$\begin{pmatrix} E_S^* b(x) E_S & A^* \\ A & -B \alpha^{-1} B^* \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix} \quad (19)$$

where $b(x)$ is either $1 + \gamma \chi_{y < 0}(x)$ for the penalty method, or $1 + l''(y(x); \mu)$ for the barrier method.

Also this problem will fit into our theoretical framework and (19) can be written in the form (5), if we set $M := \sqrt{b} E_S$ and $C^* := \alpha^{-1/2} B^*$. Clearly, also combinations of control and state constraints can be treated.

3 Two strategies for operator preconditioning

In the following we shall derive and justify some operator preconditioners for solving the minimization problem (6) with the preconditioned conjugate gradient method, see 2.3. As explained in 2.4 the preconditioners can also be used within a preconditioned MINRES method.

3.1 Preconditioning via the pure differential operators

Our first preconditioner, which is similar to the ones proposed in [10], is defined via the first part of $\langle \cdot, \cdot \rangle_K$:

$$(Q_0 v)(w) := \langle v, w \rangle_{Q_0} := \langle M^{-*} A^* v, M^{-*} A^* w \rangle_H \quad \forall w \in D_K. \quad (20)$$

If $\ell \in P^*$, then its inverse can be applied to ℓ by computing $Q_0^{-1} \ell = A^{-*} M^* M A^{-1} \ell$. However, since in general $\ell \in D_K^*$, we have to use the extended operators from Section 2.2 to make this rigorous:

$$v := Q_0^{-1} \ell = A^{-*} M^* \widehat{M} \widehat{A}^{-1} \ell. \quad (21)$$

Application of \widehat{A}^{-1} involves one solve of the extended state equation, and to apply A^{-*} one has to solve the adjoint equation. Clearly we have $A^*v = M^*\widehat{M}\widehat{A}^{-1}\ell \in \text{ran } M^*$, hence $v \in D_K$.

Lemma 3.1. *Assume that*

$$\gamma_{Q_0} := \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U}{\langle v, v \rangle_{Q_0}} \quad (22)$$

is finite. Then

$$\langle v, v \rangle_{Q_0} \leq \langle v, v \rangle_K \leq (1 + \gamma_{Q_0})\langle v, v \rangle_{Q_0}, \quad (23)$$

and thus the condition number κ_{Q_0} of K , relative to Q_0 is bounded by

$$\kappa_{Q_0} \leq (1 + \gamma_{Q_0}).$$

Proof. The left part of (23) follows simply from the positive semi-definiteness of $\langle C^*\cdot, C^*\cdot \rangle_U$. Further, since $\langle C^*v, C^*v \rangle_U \leq \gamma_{Q_0}\langle v, v \rangle_{Q_0}$, we obtain the right part. \square

3.2 A preconditioner with increased robustness

Our next preconditioner exploits the positive definiteness of CC^* to improve our condition number estimate. In order to render it well defined, we have to impose the following assumption, which restricts the class of problems to be tackled:

Assumption 3.2. (Compatibility Assumption) Let I be a non-zero continuous mapping

$$I : H \rightarrow U.$$

with $\|I\| \leq 1$.

In the case $I = 0$ the preconditioner Q that we will define below coincides with the simple choice Q_0 , defined above. By excluding this, we make sure that $Q \neq Q_0$.

For the Hilbert space adjoint $I^* : U \rightarrow H$, defined by

$$\langle I^*u, h \rangle_H = \langle u, Ih \rangle_U$$

we note that $\|I^*\| = \|I\| \leq 1$.

In all our applications I is defined as in the following example:

Example 3.3. Let Ω_H and Ω_U be two subsets of \mathbb{R}^d of non-zero measure. Define $H = L_2(\Omega_H)$ and $U = L_2(\Omega_U)$. Then the mapping $I : H \rightarrow U$ can be defined by restriction of $h \in H$ to $\Omega_H \cap \Omega_U$, followed by extension by zero onto Ω_U . In turn, $I^* : U \rightarrow H$ is the restriction to $\Omega_H \cap \Omega_U$ and extension by zero to Ω_H , namely for all $h \in L_2(\Omega_H), u \in L_2(\Omega_U)$

$$\langle u, Ih \rangle_{L_2(\Omega_U)} = \int_{\Omega_U} u(Ih) dx = \int_{\Omega_H \cap \Omega_U} uh dx = \int_{\Omega_H} (I^*u)h dx = \langle I^*u, h \rangle_{L_2(\Omega_H)}.$$

The extreme cases are following: if $\Omega_H \cap \Omega_U = \emptyset$, then $I \equiv 0$ and Assumption 3.2 is violated. The other, most desirable, extreme case is $\Omega_H = \Omega_U$ so that $I = id$.

By our assumptions the composition $CIM : Y \rightarrow P^*$ is well defined, and thus also its adjoint $(CIM)^* : P \rightarrow Y^*$ via the relations

$$(CIMy)(p) = \langle IMy, C^*p \rangle_U = \langle My, I^*C^*p \rangle_H = (M^*I^*C^*p)(y) = ((CIM)^*p)(y).$$

Definition 3.4. Define the preconditioner $Q : D_K \rightarrow D_K^*$ by:

$$(Qv)(w) := \langle v, w \rangle_Q := \langle M^{-*}(A + CIM)^*v, M^{-*}(A + CIM)^*w \rangle_H \quad \forall w \in D_K. \quad (24)$$

Since, as we will show, $\langle \cdot, \cdot \rangle_Q$ and $\langle \cdot, \cdot \rangle_K$ are equivalent (with condition number κ_Q to be estimated), the inverse $Q^{-1} : D_K^* \rightarrow D_K$ needed in the cg iteration is well defined. This follows from the Lax-Milgram theorem applied on the Hilbert space $(D_K, \langle \cdot, \cdot \rangle_K)$.

The computation of $Q^{-1}\ell$ proceeds in three steps, two of which correspond to PDE solves. In the case $\ell \in P^*$ the first step would be to solve the modified PDE $(A + CIM)y = \ell$. However, since in general $\ell \in D_K^*$, this extends to the problem (cf. Section 2.2)

$$(\widehat{A} + CIM\widehat{M})y = \ell \quad (25)$$

for $y \in \widehat{Y}$. Then one has to compute

$$w := M^*\widehat{M}y. \quad (26)$$

So that $w \in \text{ran } M^*$. Finally, one has to solve

$$(A + CIM)^*p = w. \quad (27)$$

so that $Q^{-1}\ell := p$. Since $w \in \text{ran } M^*$ and $(CIM)^*p = M^*I^*C^*p \in \text{ran } M^*$ we conclude that $A^*p \in \text{ran } M^*$, hence $p \in D_K$.

An equivalent preconditioner has been proposed and analyzed recently for the unconstrained case of distributed control by [28]. For the case of regularized state constraints an equivalent preconditioner has been proposed recently and independently in [17], but no useful estimates for the condition number were derived.

Remark 3.5. Already at this point we can predict the main features of this type of preconditioner. In contrast to Q_0 it also includes the operator C in its formulation. Hence, more information of the problem enters into the construction of the preconditioner. We will see that this leads to a significant improvement of condition numbers, in cases where Q can be applied.

However, we also observe the main limitations of our approach. The composition CIM has to be non-zero, otherwise $Q = Q_0$. Here our main focus is restricted to simple mappings I , as defined in Example 3.3, because CIM has to be simple enough to make the involved PDEs solvable at low cost.

The following lemma plays a pivotal role in our analysis:

Lemma 3.6. *Assume that the following quantity is finite:*

$$\gamma_Q := \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U}{\langle v, v \rangle_Q}. \quad (28)$$

Then we have the following estimates:

$$\frac{1}{2}\langle v, v \rangle_Q \leq \langle v, v \rangle_K \leq (2 + 3\gamma_Q)\langle v, v \rangle_Q, \quad (29)$$

and thus the condition number κ_Q of K , relative to Q is bounded by

$$\kappa_Q \leq 4 + 6\gamma_Q.$$

Proof. For the proof we recall the parallelogram law in Hilbert spaces:

$$\|x + y\|^2 + \|x - y\|^2 = 2(\|x\|^2 + \|y\|^2),$$

which implies that each summand on the left hand side can be estimated from above by the right hand side.

We use this estimate to compute

$$\begin{aligned} \langle v, v \rangle_Q &= \|M^{-*}A^*v + I^*C^*v\|_H^2 \leq 2(\|M^{-*}A^*v\|_H^2 + \|I^*C^*v\|_H^2) \\ &\leq 2(\|M^{-*}A^*v\|_H^2 + \|C^*v\|_U^2) = 2\langle v, v \rangle_K. \end{aligned}$$

Let us consider the opposite direction. For given $x \in D_K$ define

$$v := (M^{-*}A^* + I^*C^*)x$$

so that $\langle v, v \rangle_H = \langle x, x \rangle_Q$. Then, using the definition of γ_Q , we can compute:

$$\begin{aligned} \langle x, x \rangle_K &= \|M^{-*}A^*x\|_H^2 + \|C^*x\|_U^2 = \|v - I^*C^*x\|_H^2 + \|C^*x\|_U^2 \\ &\leq 2\|v\|_H^2 + 2\|I^*C^*x\|_H^2 + \|C^*x\|_U^2 \leq 2\langle x, x \rangle_Q + 3\|C^*x\|_U^2 \\ &\leq (2 + 3\gamma_Q)\langle x, x \rangle_Q. \end{aligned} \tag{30}$$

□

In order to motivate the quantity γ_Q further, consider the auxiliary quantity r , defined by

$$(A + CIM)^*v = r,$$

which means in turn that $v(r)$ is the solution of a partial differential equation with right hand side r . Then γ_Q can be written as follows:

$$\gamma_Q = \sup_{r \in \text{ran } M^*} \frac{\langle C^*v(r), C^*v(r) \rangle_U}{\langle M^{-*}r, M^{-*}r \rangle_H}.$$

Hence, our task will be to establish estimates of the form

$$\|C^*v\|_U \leq c(A, CIM)\|M^{-*}r\|_H$$

on the solution $v(r)$ of the above PDE in terms of r . This reduces the estimation of the condition number of our saddle-point system to an estimate for a PDE solution. We will use this technique in Sections 5 and 6 below, where the particular structure of the PDE at hand is used. In the following and in Section 4, we keep arguing purely in terms of functional analytic estimates.

Sharpness of Lemma 3.6. With some additional effort one can refine the estimate (29) slightly. However, for us, the asymptotics $\kappa_Q = O(\gamma_Q)$ for large γ_Q is the main point of interest. This relation cannot be improved substantially, as we will briefly explain.

Consider the left bound first. Since usually A^* is a differential operator, we may assume that there is a sequence v_k , such that $\|M^{-*}A^*v_k\|_H = 1$, while $\|C^*v_k\|_U \rightarrow 0$. Then the terms in K and Q containing C^* can be neglected, and we obtain $\langle v_k, v_k \rangle_Q / \langle v_k, v_k \rangle_K \rightarrow 1$.

As for the right bound in (29), v can be chosen such that $(1 - \varepsilon)\gamma_Q$ is attained in (28) for sufficiently small ε . Then instead of (30) we can compute

$$\langle v, v \rangle_K = \|M^{-*}A^*v\|_H^2 + \|C^*v\|_U^2 \geq \|C^*v\|_U^2 \geq (1 - \varepsilon)\gamma_Q \langle v, v \rangle_Q.$$

Hence, taking both estimates together we obtain a lower bound for the condition number, given by

$$\kappa_Q \geq \gamma_Q. \tag{31}$$

4 Applications to concrete problems

In this section we will discuss a couple of examples for which our preconditioning strategy can (or cannot) be applied effectively. This should clarify the advantages, but also the limitations of the preconditioner Q from (24). Included are problems with control constraints, and also with state constraints. The following bounds hold in a quite general setting. Under stronger assumptions they can be refined, as shown in Section 5 and Section 6.

We will proceed in two steps. First we will exploit the properties of the differential operator A , corresponding to elliptic or parabolic PDEs to give a more concrete estimate for γ_Q . Second, we will use these estimates to compute condition numbers for a selection of examples.

4.1 Estimating γ_Q for elliptic and parabolic problems

In this section we will establish estimates on γ_Q , exploiting the properties of the differential operators A and A^* .

We will show that both elliptic and parabolic problems admit estimates of the common form

$$\gamma_Q \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle \langle Mv, Mv \rangle \rangle}{(\langle \langle v, v \rangle \rangle + \langle \langle C^*v, IMv \rangle \rangle)^2}, \quad (32)$$

where the notation $\langle \langle \cdot, \cdot \rangle \rangle$ stands for one of the scalar products $\langle \cdot, \cdot \rangle_H$ and $\langle \cdot, \cdot \rangle_{e^{-\omega t}}$ (defined below), which are suited for the elliptic and the parabolic case, respectively.

In fact, (32) follows from (33) and (35) by taking into account that A^* is positive definite in the corresponding scalar product, i.e.,

$$(A^*v)(v) \geq c_A \langle \langle v, v \rangle \rangle.$$

In the elliptic case, this is true due to ellipticity of $(A^*v)(v) = a(v, v)$, in the parabolic case, this follows from (36), below.

4.1.1 The elliptic case

Consider the preconditioner Q from (24) for the elliptic case, where $A : H^1(\Omega) \rightarrow H^1(\Omega)^*$ is defined via an elliptic bilinear form $a(\cdot, \cdot)$ as in (17) from subsection 2.5.1.

Lemma 4.1. *For the elliptic equation (17) we have the estimate*

$$\gamma_Q \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_H}{((A^*v)(v) + \langle C^*v, IMv \rangle_U)^2} \quad (33)$$

Proof. Let $v \in D_K \subset P = H^1(\Omega)$, so that $(A^* + M^*I^*C^*)v \in \text{ran } M^* \subset Y^*$. We start with the Cauchy-Schwarz inequality:

$$\begin{aligned} ((A^* + M^*I^*C^*)v)(v) &= \langle M^{-*}(A^* + M^*I^*C^*)v, Mv \rangle_H \\ &\leq \|M^{-*}(A^* + M^*I^*C^*)v\|_H \|Mv\|_H \\ &= \sqrt{\langle v, v \rangle_Q} \langle Mv, Mv \rangle_H. \end{aligned}$$

Hence, we may estimate

$$\begin{aligned} \gamma_Q &= \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_H}{\langle v, v \rangle_Q \langle Mv, Mv \rangle_H} \leq \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_H}{((A^* + M^*I^*C^*)v)(v)^2} \\ &= \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_H}{((A^*v)(v) + \langle C^*v, IMv \rangle_U)^2}. \end{aligned}$$

□

4.1.2 The parabolic case

To establish an estimate for γ_Q in the parabolic case, as in subsection 2.5.2, we have to modify our proof slightly in order to cope with the non-symmetry of its differential operator, which reads

$$(Ay)(v) = \int_{[0,T]} \left(\langle y_t, v \rangle + \int_{\Omega} \langle \nabla y, \sigma(t, x) \nabla v \rangle_{\mathbb{R}^d} + a_0(t, x) y v \, dx \right) dt, \quad (34)$$

together with initial conditions $y(0) = 0$. We have to employ a special scalar and duality product. For $\omega > 0$ we set

$$\langle v, w \rangle_{e^{-\omega t}} := \int_{[0,T]} e^{-\omega t} \langle v(t), w(t) \rangle_{L_2(\Omega)} dt,$$

which induces an equivalent norm $e^{-\omega T} \|\cdot\|_{L_2([0,T] \times \Omega)} \leq \|\cdot\|_{e^{-\omega t}} \leq \|\cdot\|_{L_2([0,T] \times \Omega)}$. Similarly, for a Banach space V we write the duality product on $L_2([0, T], V^*) \times L_2([0, T], V)$

$$\langle v^*, v \rangle_{e^{-\omega t}} := \int_{[0,T]} e^{-\omega t} (v^*(t)(v(t))) dt,$$

Our motivation is that A^* is positive definite with respect to this scalar product, as long as ω is chosen sufficiently large, as will be shown in the first part of the proof of the next result (in fact this is a standard result in the theory of parabolic equations).

Lemma 4.2. *Let A be defined as in (34). Assume that $\langle M^* v, w \rangle_{e^{-\omega t}} = \langle v, M w \rangle_{e^{-\omega t}}$. Then A^* is positive definite w.r.t. $\langle \cdot, \cdot \rangle_{e^{-\omega t}}$, and we obtain the following condition number:*

$$\gamma_Q \leq c(T) \sup_{v \in D_K} \frac{\langle C^* v, C^* v \rangle_U \langle M v, M v \rangle_{e^{-\omega t}}}{(\langle A^* v, v \rangle_{e^{-\omega t}} + \langle I^* C^* v, M v \rangle_{e^{-\omega t}})^2} \quad (35)$$

Proof. First, we show that A^* is positive definite w.r.t. the scalar product $\langle \cdot, \cdot \rangle_{e^{-\omega t}}$. Inserting $w = e^{-\omega t} v$ into the formula of integration by parts (cf. e.g. [8, Satz 1.17])

$$\langle v(T), w(T) \rangle - \langle v(0), w(0) \rangle = \int_{[0,T]} \langle v_t(t), w(t) \rangle + \langle w_t(t), v(t) \rangle dt$$

and taking into account our restriction $v(0) = w(0) = 0$ we infer after a short computation

$$\langle v_t, v \rangle_{e^{-\omega t}} = \frac{1}{2} \left(e^{-\omega T} \|v(T)\|^2 + \omega \langle v, v \rangle_{e^{-\omega t}} \right).$$

As for the remaining part of A^* we have

$$\int_{[0,T]} \int_{\Omega} \langle \nabla v, \sigma(t, x) \nabla v \rangle_{\mathbb{R}^d} + a(t, x) v^2 \, dx \, e^{-\omega t} dt \geq 0,$$

and hence

$$\langle A^* v, v \rangle_{e^{-\omega t}} \geq \frac{\omega}{2} \langle v, v \rangle_{e^{-\omega t}}. \quad (36)$$

From this point our proof runs in parallel to the elliptic case. Similar as before, let $v \in Y = W([0, T])$. Then also $v \in P = L_2([0, T], H^1(\Omega))$, such that $(A^* + M^*I^*C^*)v \in Y^*$. Thus, we can use the Cauchy-Schwarz inequality:

$$\begin{aligned} & \langle (A^* + M^*I^*C^*)v, v \rangle_{e^{-\omega t}}^2 \\ & \leq \langle M^{-*}(A^* + M^*I^*C^*)v, M^{-*}(A^* + M^*I^*C^*)v \rangle_{e^{-\omega t}} \langle Mv, Mv \rangle_{e^{-\omega t}} \\ & \leq \langle v, v \rangle_Q \langle Mv, Mv \rangle_{e^{-\omega t}}. \end{aligned}$$

Hence, we may estimate, also as before:

$$\begin{aligned} \gamma_Q &= \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_{e^{-\omega t}}}{\langle v, v \rangle_Q \langle Mv, Mv \rangle_{e^{-\omega t}}} \leq \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_{e^{-\omega t}}}{\langle (A^* + M^*I^*C^*)v, v \rangle_{e^{-\omega t}}^2} \\ &= \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle Mv, Mv \rangle_{e^{-\omega t}}}{(\langle A^*v, v \rangle_{e^{-\omega t}} + \langle I^*C^*v, Mv \rangle_{e^{-\omega t}})^2}. \end{aligned}$$

□

Remark 4.3. Recall that $M^* : H \rightarrow Y^*$ is the adjoint of $M : Y \rightarrow H$ w.r.t. $\langle \cdot, \cdot \rangle_H$, so our assumption $\langle M^*v, w \rangle_{e^{-\omega t}} = \langle v, Mw \rangle_{e^{-\omega t}}$ is needed. However, it can easily be verified, if for example M can be written as $(Mv)(t) = M(t)v(t)$, where $M(t)$ depends on the “slice” $v(t)$ only.

The dependence of γ_Q on the interval length T can be worked out to be proportional to T by choosing ω optimally.

4.2 Application to concrete examples

Let us now continue the discussion of the examples, presented in Section 2.5.

4.2.1 Distributed control problems with control bounds

Let us consider as an example the optimal control problem

$$\min \frac{1}{2} \|My - y_d\|_{L_2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L_2(\Omega)}^2 \quad \text{s.t. } Ay - Bu = 0 \quad u \geq 0$$

As explained in Section 2.5.1 this problem can be solved by a semi-smooth Newton method, which leads to systems of the form (4). Recall that in this example $Y = P = H_0^1(\Omega)$ and $U = H = L_2(\Omega)$ as well as the definition of $a(\cdot, \cdot)$ as an elliptic bilinear form on $H_0^1(\Omega)$. Further, $M : H_0^1(\Omega) \rightarrow L_2(\Omega)$ is the Sobolev embedding E_S , $B = E_S^* : L_2(\Omega) \rightarrow H_0^1(\Omega)^*$, and $C^* : H_0^1(\Omega) \rightarrow L_2(\Omega)$ is defined as $C^* = \alpha^{-1/2} \chi_{\mathcal{I}}(p) E_S$. Finally, the mapping $I : H \rightarrow U$ is just the identity in $L_2(\Omega)$.

Taking this into account, we can write down the concrete representation

$$((A + CIM)y)(p) = ((A + CIM)^*p)(y) = a(y, p) + \int_{\Omega} \phi y p \, dx \quad \text{for } y, p \in H_0^1(\Omega), \quad (37)$$

where $\phi(x) = \alpha^{-1/2} \chi_{\mathcal{I}}(x) \in L_{\infty}(\Omega)$. Hence, the steps (25) and (27) during the application of the preconditioner amount in solving modified PDEs with an additional mass-term, defined by ϕ .

To summarize, the computation of $Q^{-1}\ell$ requires the solution of two Poisson problems, the first of which is in very weak form, the second in strong form.

With our analysis from the previous section we obtain the following results for our preconditioners.

Proposition 4.4. *Consider the preconditioner Q_0 from (20) applied to the block operator (4). Then we obtain the following condition number:*

$$\kappa_{Q_0} \leq 1 + c\alpha^{-1}.$$

Proof. In view of Lemma 3.1 we have to provide estimates for γ_{Q_0} in (22). For the numerator we can compute

$$\langle C^*v, C^*v \rangle_U = \langle \alpha^{-1}\chi_{\mathcal{I}}(p)v, v \rangle_{L_2} \leq \|\alpha^{-1}\chi_{\mathcal{I}}(p)\|_{\infty} \|v\|_{L_2}^2$$

and the denominator yields $\langle M^{-*}A^*v, M^{-*}A^*v \rangle_H \geq c\|v\|_{L_2}^2$ by continuity of $(A^*)^{-1} : L_2(\Omega) \rightarrow L_2(\Omega)$. So $\gamma_{Q_0} \leq c\alpha^{-1}$, which yields the desired result for κ_{Q_0} via Lemma 3.1. \square

Now we consider our preconditioner Q for the elliptic and parabolic case.

Proposition 4.5. *Consider the preconditioner Q from (24) applied to the block operator (4). Consider the elliptic or the parabolic operator A from Section 4.1. Then*

$$\kappa_Q \leq c(1 + \alpha^{-1/2})$$

Proof. As in (32) let $\langle\langle \cdot, \cdot \rangle\rangle$ be one of $\langle \cdot, \cdot \rangle_{L_2(\Omega)}$ (for the elliptic case) or $\langle \cdot, \cdot \rangle_{e^{-\omega t}}$ (for the parabolic case). Then, from (32) we estimate (taking into account $Mv = v$).

$$\gamma_Q \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle\langle v, v \rangle\rangle}{(\langle\langle v, v \rangle\rangle + \langle\langle C^*v, v \rangle\rangle)^2} \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle\langle v, v \rangle\rangle}{2\langle\langle v, v \rangle\rangle \langle\langle C^*v, v \rangle\rangle}.$$

The last inequality follows from the general relation $a^2 + b^2 \geq 2ab$. Moreover, in both cases, $\langle C^*v, C^*v \rangle_U \leq c\langle\langle C^*v, C^*v \rangle\rangle$. By definition of C we obtain

$$\langle C^*v, C^*v \rangle_U \leq c\langle\langle C^*v, C^*v \rangle\rangle \leq c\alpha^{-1/2}\langle\langle C^*v, v \rangle\rangle.$$

Hence,

$$\gamma_Q \leq c\alpha^{-1/2},$$

and thus by Lemma 3.6 $\kappa_Q \leq c(1 + \alpha^{-1/2})$. \square

Thus, we have obtained $\kappa_Q \sim \sqrt{\kappa_{Q_0}}$ which already yields a considerable gain of efficiency for a cg method via (12).

If no bounds on the control are present we can recover the results obtained in [28]:

Corollary 4.6. *In the unconstrained case we obtain the α independent bound*

$$\kappa_Q \leq c$$

Proof. Just as before, we compute

$$\gamma_Q = \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_U \langle\langle v, v \rangle\rangle}{(c\langle\langle v, v \rangle\rangle + \langle\langle C^*v, v \rangle\rangle)^2} \leq c \sup_{v \in D_K} \frac{\alpha^{-1}\langle\langle v, v \rangle\rangle^2}{\langle\langle C^*v, v \rangle\rangle^2} = c \sup_{v \in D_K} \frac{\alpha^{-1}\langle\langle v, v \rangle\rangle^2}{(\alpha^{-1/2}\langle\langle v, v \rangle\rangle)^2}.$$

Hence $\gamma_Q \leq c$. \square

The parabolic case. As a parabolic example we consider the one from Section 2.5.2:

$$\min \frac{1}{2} \|My - y_d\|_{L_2([0,T], L_2(\Omega))}^2 + \frac{\alpha}{2} \|u\|_{L_2([0,T], L_2(\Omega))}^2 \text{ s.t. } Ay - Bu = 0 \quad u \geq 0.$$

Here, $Y = W([0, T])$, $P = L_2([0, T], H^1)$ and $U = H = L_2([0, T], L_2(\Omega))$. Similar to before, $M : Y \rightarrow H$ is the Sobolev embedding $W([0, T]) \hookrightarrow L_2([0, T], L_2(\Omega))$. Using the Sobolev embedding $E_S : L_2([0, T], H^1(\Omega)) \rightarrow L_2([0, T], L_2(\Omega))$ we can define $B = E_S^* : U \rightarrow P^*$, and $C^* : P \rightarrow U$ as $C^* = \alpha^{-1/2} \chi_{\mathcal{I}}(p) E_S$. Also here, the mapping $I : H \rightarrow U$ is just the identity on $L_2([0, T], L_2(\Omega))$.

In the parabolic case, the application of the preconditioner comprises a forward and a backward solve of modified parabolic PDEs. By the same argumentation as in the elliptic case (25) can be extended to right-hand sides in D_K^* .

4.2.2 Disjoint control and observation regions

Consider the problem (for simplicity, let A be the elliptic operator from (17))

$$\min \frac{1}{2} \|y - y_d\|_{L_2(\Omega_H)}^2 + \frac{\alpha}{2} \|u\|_{L_2(\Omega_U)}^2 \text{ s.t. } Ay - Bu = 0,$$

where $B : L_2(\Omega_U) \rightarrow P^*$ is a continuous mapping, and Ω_H and Ω_U are subsets of $\bar{\Omega}$, such that their intersection is a set of measure zero in both $H := L_2(\Omega_H)$ and $U := L_2(\Omega_U)$.

An important special case is boundary control and observation in the domain, i.e., $\Omega_H = \Omega$ and $\Omega_U = \Gamma = \partial\Omega$. The only simple choice is $I : L_2(\Omega) \rightarrow L_2(\Gamma)$ via $I \equiv 0$, which violates Assumption 3.2. Thus, only Q_0 can be used. A similar case is boundary observation and distributed control.

4.2.3 Distributed control with state constraints

In the case of state constrained optimal control problems, algorithms often employ a path-following scheme, which leads to a block M that is very ill conditioned towards the end of the algorithm. Both our preconditioners can be applied in this case. Moreover, it can be shown that usually the preconditioner Q is significantly more robust than Q_0 with respect to the path-following parameters.

As elaborated in Section 2.5.4 in state constrained problems with distributed control linear systems of the form

$$\begin{pmatrix} E_S^* b(x) E_S & A^* \\ A & -B\alpha^{-1} B^* \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix} \quad (38)$$

have to be solved, where $E_S = B^*$ is the Sobolev embedding, in the case of distributed control, where $b(x) > 0$, and typically $\|b\|_\infty$ is very large, tending to infinity towards the end of the algorithm. Just as in Section 4.2.1 we can define $H = U = L_2(\Omega)$ and $M = \sqrt{b(x)} E$ and $C^* = \alpha^{-1/2} B$ and again $I = Id$. Also the Definition of D_K and the application on Q^{-1} via extensions of (25)-(27) are the same as in Section 4.2.1.

Only the condition number estimates are different:

Proposition 4.7. *Assume that $p \in L_\infty$. Consider the preconditioner Q_0 from (20) applied to the operator (38). Then we obtain the following condition number:*

$$\kappa_{Q_0} \leq 1 + c \|b\|_\infty \alpha^{-1}$$

Proof. By Lemma 3.1 we have to provide estimates for γ_{Q_0} in (22). For the numerator we can compute

$$\langle C^*v, C^*v \rangle \leq \alpha^{-1} \|v\|_{L_2}^2$$

and the denominator yields $\langle M^{-*}A^*v, M^{-*}A^*v \rangle \geq c \|b\|_\infty^{-1} \|v\|_{L_2(\Omega)}^2$ by continuity of $(A^*)^{-1} : L_2(\Omega) \rightarrow L_2(\Omega)$. \square

The following result crucially depends on the assumption that C and M are defined as multiplication operators with functions that have the same support. This is true in particular for purely state constrained problems with distributed control.

Proposition 4.8. *Consider the preconditioner Q from (24) applied to the block operator (38). Then the condition number is bounded by*

$$\kappa_Q \leq c(1 + \|b\|_\infty^{1/2} \alpha^{-1/2})$$

Proof. Inserting $\langle \cdot, \cdot \rangle_U = \langle \cdot, \cdot \rangle_H = \langle \cdot, \cdot \rangle_{L_2}$ and $(A^*v)(v) \geq \langle v, v \rangle_{L_2}$ in Lemma 4.1 yields

$$\gamma_Q \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_{L_2} \langle Mv, Mv \rangle_{L_2}}{(\langle v, v \rangle_{L_2} + \langle C^*v, IMv \rangle_{L_2})^2} \leq c \sup_{v \in D_K} \frac{\langle C^*v, C^*v \rangle_{L_2} \langle Mv, Mv \rangle_{L_2}}{2 \langle v, v \rangle_{L_2} \langle C^*v, IMv \rangle_{L_2}}. \quad (39)$$

By definition of C and M we obtain

$$\begin{aligned} \langle C^*v, C^*v \rangle_{L_2} &= \alpha^{-1} \langle v, v \rangle_{L_2} \\ \langle Mv, Mv \rangle_{L_2} &\leq \alpha^{1/2} \|b\|_\infty^{1/2} \langle C^*v, IMv \rangle_{L_2}. \end{aligned}$$

Here we used that $M^*I^*C^*v = \alpha^{-1/2}(b(x))^{1/2}v$ is just a pointwise multiplication by positive functions. Hence, we finally compute

$$\gamma_Q \leq c\alpha^{-1/2} \|b\|_\infty^{1/2}.$$

\square

A similar situation holds for boundary control problems if the state constraints are only imposed on the boundary. For parabolic problems, a similar result is obtained analogously, replacing Lemma 4.1 by Lemma 4.2.

The improved preconditioner Q is not effective for boundary control and state constraints, unless these are also imposed on the boundary only. In that case, we have again $CIM = 0$, and thus $Q_0 = Q$ just as described in Section 4.2.2. As a remedy, it is conceivable to introduce an artificial “virtual” control [15] on the domain equipped with a regularization parameter that is driven to ∞ .

A similar situation occurs for additional control constraints, in case that the control is active, i.e. $\chi_{\mathcal{I}}(p)(x) = 0$ where $b(x)$ is large. It is, however, still feasible in this case to use Q as a preconditioner. If active control and state set are disjoint, then we conjecture that Q is still more efficient than Q_0 . Otherwise, its efficiency may degrade to a level comparable to Q_0 if this assumption is not valid.

5 A splitting argument for refined estimates

In some situations our estimates for Q can still be refined, if we are willing to impose additional assumptions on CIM . In broad terms, we will assume that CIM is defined as multiplication operators via piecewise constant functions with smoothly bounded level

sets. The applications below comprise Newton systems for control constrained problems and for penalty methods for state constrained problems as described in Example 2.5.1 and Example 2.5.4. For simplicity we concentrate on the elliptic case.

In the following, let $a : H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$ be an elliptic bilinear form, as defined, e.g., in (17):

$$a(y, v) = \int_{\Omega} \langle \sigma(x) \nabla y, \nabla v \rangle_{\mathbb{R}^a} + a_0(x) y v \, dx \quad (40)$$

Further, let us define the corresponding operator in strong form:

$$\mathcal{A}y := -\operatorname{div}(\sigma(x) \nabla y) + a_0(x) y, \quad (41)$$

which results from $a(\cdot, \cdot)$ via integration by parts, assuming that $\sigma(x)$ is sufficiently smooth.

Assumption 5.1. Assume that CIM is a multiplication operator, defined by a piecewise constant function ϕ , i.e., $(CIMv)(x) = \phi(x)v(x)$, which assumes two non-negative values ϕ_1 and ϕ_0 , such that $\phi_1 > 0$ and $\phi_1 \geq \phi_0 \geq 0$, namely:

$$\phi(x) = \begin{cases} \phi_0 & : x \in J_0 \\ \phi_1 & : x \in J_1 := \Omega \setminus J_0. \end{cases} \quad (42)$$

Let us denote by ∂J_0 the boundary of J_0 relative to Ω , i.e., $\partial J_0 \subset \Omega$, so that $\partial J_0 = \partial J_1$.

Assume that J_0 and J_1 are Lipschitz domains and that the solution v_J of the problem

$$v_J \in H_0^1(J_0) : \quad a(v_J, w) = \int_{J_0} f w \, dx \quad \forall w \in H_0^1(J_0)$$

gives rise to the following trace estimate:

$$\|\partial_{\sigma\nu} v_J\|_{L_2(\partial J_0)} \leq c_{tr,1} \|f\|_{L_2(J_0)}. \quad (43)$$

Here $\partial_{\sigma\nu}$ stands for the derivative in direction of the outer normal ν at a point $x \in \partial J_0$ with respect to the scalar product induced by $\sigma(x)$. This normal has to be defined almost everywhere on ∂J_0 .

The validity of (43) certainly depends on the smoothness of ∂J_0 and on the coefficients of $a(\cdot, \cdot)$, and is known to hold e.g., for $H^{3/2+\varepsilon}$ -regular problems. In general one only has (43) for $\|\partial_{\sigma\nu} v_J\|_{H^{-1/2}(\partial J_0)}$. In the context of optimal control the smoothness of ∂J_0 can usually only be observed a-posteriori, for example, after the active set of the optimal control has been computed. From an intuitive point of view, smoothness of ∂J_0 results in a relatively weak coupling of J_0 and J_1 . So it should have a certain influence on the condition number.

Lemma 5.2. *Let J be an open domain with Lipschitz boundary. Then for $v \in H^1(J)$ the following estimate holds for the trace operator $\tau : H^1(J) \rightarrow L_2(\partial J)$:*

$$\|\tau(v)\|_{L_2(\partial J)} \leq c_{tr,2} \sqrt{\|v\|_{H^1(J)} \|v\|_{L_2(J)}}.$$

Proof. Since $v \in H^1(J)$, $\tau(v)$ exists by the classical trace theorem. After localization and transformation of a part of the boundary to the first coordinate axis, we end up in showing (dividing the coordinates into $x = (t, x')$):

$$\|v(0, \cdot)\|_{L_2}^2 \leq c \|v\|_{H^1(J)} \|v\|_{L_2(J)}.$$

This can be obtained by the formula of integration by parts, well known from parabolic problems:

$$2 \int_{[0,T]} \left\langle \frac{d}{dt} v(t), v(t) \right\rangle dt = \|v(T, \cdot)\|_{L_2}^2 - \|v(0, \cdot)\|_{L_2}^2$$

Choosing T large enough, such that $v(T) \equiv 0$, we obtain

$$\|v(0, \cdot)\|_{L_2}^2 \leq 2 \int_{[0,T]} \left\| \frac{d}{dt} v(t) \right\|_{L_2} \|v(t)\|_{L_2} dt \leq 2 \|v\|_{H^1} \|v\|_{L_2}.$$

□

This estimate can *not* be acquired directly via interpolation theory of Sobolev spaces, because there exists no continuous trace operator $H^{1/2}(J) \rightarrow L_2(\partial J)$.

In the following lemma we will consider the problem

$$v \in H_0^1(\Omega) : \quad a(v, w) + \int_{\Omega} \phi(x) v(x) w(x) dx = \int_{\Omega} f w dx \quad \forall w \in H_0^1(\Omega). \quad (44)$$

This can also be written in operator notation, if we define (E_S is the Sobolev embedding):

$$\begin{aligned} \Phi : H_0^1(\Omega) &\rightarrow H_0^1(\Omega)^* \\ (\Phi v)(x) &:= E_S^* \phi(x) (E_S v)(x) \end{aligned} \quad (45)$$

Then (44) reads

$$(A + \Phi)v = f.$$

Lemma 5.3. *Consider problem (44) such that ϕ has the properties, defined in Assumption 5.1, and assume that $f \in L_2(\Omega)$. Then*

$$\|\Phi v\|_{L_2(\Omega)} \leq c(\|f\|_{L_2(\Omega)} + \phi_1^{1/4} \|f\|_{L_2(J_0)}). \quad (46)$$

If additionally $\phi_0 > 0$, then

$$\|v\|_{L_2(\Omega)} \leq c(1 + \phi_1^{1/4}) \|\Phi^{-1} f\|_{L_2(\Omega)}. \quad (47)$$

In all these estimates c depends on the regularity of J_0 and J_1 , and on the ellipticity of $a(\cdot, \cdot)$ in (44).

Proof. The idea of the proof is to split v into two parts $v = v_0 + v_1$. To define v_0 , we first consider v_J , the solution of the following problem on J_0 :

$$v_J \in H_0^1(J_0) : \quad a(v_J, w) + \int_{J_0} \phi_0 v_J w dx = \int_{J_0} f w dx \quad \forall w \in H_0^1(J_0). \quad (48)$$

Now we extend v_J by 0 to a function $v_0 \in H_0^1(\Omega)$, such that $v_0|_{J_0} = v_J$, and $v_0 = 0$ on J_1 . This implies also $\phi_0 v_0(x) = \phi(x) v_0(x)$ for almost all $x \in \Omega$, because $\phi(x) = \phi_0$ on J_0 and $v_0 = 0$ on J_1 .

Testing (48) with v_J and dividing by $\|v_J\|_{L_2(J_0)}$ we obtain

$$\|\phi v_0\|_{L_2(\Omega)} = \phi_0 \|v_J\|_{L_2(J_0)} \leq \|f\|_{L_2(J_0)}. \quad (49)$$

Thus, v_J satisfies

$$a(v_J, w) = \int_{J_0} \underbrace{(f - \phi_0 v_J)}_{\tilde{f}} w dx \quad \forall w \in H_0^1(J_0),$$

with $\|\tilde{f}\|_{L_2} \leq 2\|f\|_{L_2}$. By our trace assumption we conclude

$$\|\partial_{\sigma\nu}v_J\|_{L_2(\partial J_0)} \leq c_{\text{tr},1}2\|f\|_{L_2(J_0)}. \quad (50)$$

Integration by parts on J_0 yields

$$\begin{aligned} \mathcal{A}v_0 + \phi_0v_0 &= f \text{ on } J_0 \\ v_0 &= 0 \text{ on } \partial J_0 \cup \partial\Omega \cup J_1. \end{aligned}$$

Testing this equation with $w \in H_0^1(\Omega)$ and separate integration by parts on J_0 and J_1 reveals that v_0 (recall $\phi v_0 = \phi_0v_0$ satisfies the weak form

$$v_0 \in H_0^1(\Omega) : a(v_0, w) + \int_{J_0} \phi(x)v_0w \, dx + \int_{\partial J_0} \partial_{\sigma\nu}v_Jw \, ds = \int_{J_0} fw \, dx \quad \forall w \in H_0^1(\Omega). \quad (51)$$

Hence, if we define v_1 as the solution of the following problem:

$$v_1 \in H_0^1(\Omega) : a(v_1, w) + \int_{\Omega} \phi(x)v_1w \, dx - \int_{\partial J_0} \partial_{\sigma\nu}v_Jw \, ds = \int_{J_1} fw \, dx \quad \forall w \in H_0^1(\Omega), \quad (52)$$

we see (by adding (51) and (52)) that $v = v_0 + v_1$ solves our original problem (44).

To obtain an estimate for v_1 we test (52) with v_1 and get

$$a(v_1, v_1) + \int_{\Omega} \phi(x)v_1^2 \, dx \leq \|\partial_{\sigma\nu}v_J\|_{L_2(\partial J_0)}\|v_1\|_{L_2(\partial J_0)} + \|f\|_{L_2(J_1)}\|v_1\|_{L_2(J_1)}. \quad (53)$$

By Lemma 5.2, we obtain:

$$\|v_1\|_{L_2(\partial J_0)} \leq c_{\text{tr},2}\sqrt{\|v_1\|_{L_2(J_1)}\|v_1\|_{H^1(J_1)}}.$$

Division of (53) by the square-root of its left-hand-side and taking into account $\phi_1 > 0$ and the ellipticity condition $a(v, v) \geq c_a\|v\|_{H^1(\Omega)}^2$ we obtain due to (50):

$$\begin{aligned} \|\phi^{1/2}v_1\|_{L_2(\Omega)} &\leq \sqrt{a(v_1, v_1) + \int_{\Omega} \phi(x)v_1^2 \, dx} \\ &\leq \frac{\|\partial_{\sigma\nu}v_J\|_{L_2(\partial J_0)}c_{\text{tr},2}\sqrt{\|v_1\|_{L_2(J_1)}\|v_1\|_{H^1(J_1)}} + \|f\|_{L_2(J_1)}\|v_1\|_{L_2(J_1)}}{\sqrt{c_a\|v_1\|_{H^1(\Omega)}^2 + \phi_1\|v_1\|_{L_2(J_1)}^2}} \\ &\leq c\|\partial_{\sigma\nu}v_J\|_{L_2(\partial J_0)}\phi_1^{-1/4} + \|f\|_{L_2(J_1)}\phi_1^{-1/2} \leq c\|f\|_{L_2(J_0)}\phi_1^{-1/4} + \|f\|_{L_2(J_1)}\phi_1^{-1/2}. \end{aligned}$$

Taking into account that $\phi_1 = \|\phi\|_{\infty}$ we obtain from this the estimate

$$\|\phi v_1\|_{L_2(\Omega)} \leq \phi_1^{1/2}\|\phi^{1/2}v_1\|_{L_2(\Omega)} \leq c\phi_1^{1/4}\|f\|_{L_2(J_0)} + \|f\|_{L_2(J_1)}. \quad (54)$$

By the triangle inequality we combine the estimates (49) for v_0 and (54) for v_1 :

$$\begin{aligned} \|\phi v\|_{L_2(\Omega)} &= \|\phi(v_0 + v_1)\|_{L_2(\Omega)} \leq \|\phi v_0\|_{L_2(\Omega)} + \|\phi v_1\|_{L_2(\Omega)} \\ &\leq \|f\|_{L_2(J_0)} + c\phi_1^{1/4}\|f\|_{L_2(J_0)} + \|f\|_{L_2(J_1)} \leq \sqrt{2}\|f\|_{L_2(\Omega)} + c\phi_1^{1/4}\|f\|_{L_2(J_0)}. \end{aligned}$$

Tracing back the constant c , we notice that it depends solely on $c_{\text{tr},1}$, $c_{\text{tr},2}$, and c_a .

For the second result we apply a duality technique. For given $w \in H_0^1(\Omega)$ define $z_w := (A^* + \Phi)^{-1}w$. Since A^* has by (44) the same properties as A , we can compute

$$\begin{aligned} \|v\|_{L_2(\Omega)} &= \sup_{\|w\|_{L_2(\Omega)}=1} \langle v, w \rangle_{L_2(\Omega)} = \sup_{\|w\|_{L_2(\Omega)}=1} \langle (A + \Phi)^{-1}f, w \rangle \\ &= \sup_{\|w\|_{L_2(\Omega)}=1} \langle f, z_w \rangle = \sup_{\|w\|_{L_2(\Omega)}=1} \langle \Phi^{-1}f, \Phi z_w \rangle \\ &\leq \sup_{\|w\|_{L_2(\Omega)}=1} \|\Phi^{-1}f\|_{L_2(\Omega)} \|\Phi z_w\|_{L_2(\Omega)} \\ &\leq \sup_{\|w\|_{L_2(\Omega)}=1} \|\Phi^{-1}f\|_{L_2(\Omega)} (\sqrt{2} + c\phi_1^{1/4}) \|w\|_{L_2(\Omega)}, \end{aligned}$$

which implies our assertion, since $\|w\|_{L_2(\Omega)} = 1$. \square

Remark 5.4. If the smoothness of ∂J_0 does not admit an L_2 -estimate of the form (43) but only in a weaker norm (e.g. in $\|\partial_{\sigma\nu} v_J\|_{H^{-s}(\partial J_0)}$ for $s \in [0, 1/2]$), one can show a similar result, where $\phi_1^{1/4}$ is replaced by $\phi_1^{1/4+s/2}$.

Sharpness of Lemma 5.3. In the following we will briefly argue that the estimate (46) is sharp. Let $\Omega =]0, 2[\subset \mathbb{R}$, and choose $\phi = \phi_0 = 0$ on $]1, 2[$ and $\phi = \phi_1 = \text{const}$ on $]0, 1[$. Moreover, set $f = 0$ on $]0, 1[$ and $f = 2$ on $[1, 2[$. Consider the problem

$$-v'' + \phi v = f \text{ on }]0, 2[, \quad v'(0) = 0, \quad v(2) = 0,$$

which is by symmetry one half of a Dirichlet problem on $] - 2, 2[$. We can now proceed along the lines of our proof and split $v = \tilde{v} + v_0$, where v_0 solves the ϕ -independent problem

$$-v_0'' = 2 \text{ on }]1, 2[, \quad v_0(1) = 0, \quad v_0(2) = 0,$$

which has the solution $v_0(x) = -(x - 1.5)^2 + 0.25$ with derivative $v_0'(1) = 1$ at $x = 1$. The second part \tilde{v} satisfies the following differential equation

$$\begin{aligned} -\tilde{v}'' + \phi_1 \tilde{v} &= 0 \text{ on } [0, 1], \quad \tilde{v}'(0) = 0 \\ -\tilde{v}'' &= 0 \text{ on }]1, 2[, \quad \tilde{v}(2) = 0 \\ \tilde{v}'_-(1) &= \tilde{v}'_+(1) + v_0'(1) = \tilde{v}'_+(1) + 1, \end{aligned}$$

where $\tilde{v}'_-(1)$ and $\tilde{v}'_+(1)$ denote the left and right limit of \tilde{v}' at $x = 1$, respectively. Obviously, \tilde{v} is a linear polynomial on $]1, 2[$, so that by our boundary conditions we have $\tilde{v}'_+(1) = -\tilde{v}(1)$, and it remains to solve the following problem on $[0, 1]$:

$$-\tilde{v}'' + \phi_1 \tilde{v} = 0 \text{ on } [0, 1], \quad \tilde{v}'(0) = 0, \quad \tilde{v}'(1) = 1 - \tilde{v}(1).$$

By a classical ansatz, this equation has a solution of the form

$$\begin{aligned} \tilde{v}(x) &= a \cosh(\sqrt{\phi_1}x) \\ \tilde{v}'(x) &= \sqrt{\phi_1}a \sinh(\sqrt{\phi_1}x) \end{aligned}$$

which already has $\tilde{v}'(0) = 0$ built in, so that we only have to determine a from the condition $\tilde{v}'(1) = 1 - \tilde{v}(1)$. A short computation yields:

$$a = \frac{1}{\sqrt{\phi_1} \sinh(\sqrt{\phi_1}) + \cosh(\sqrt{\phi_1})},$$

so that we can compute (recall $\phi_0 = 0$)

$$\begin{aligned} \|\phi v\|_{L_2([0,2])}^2 &= \|\phi_1 \tilde{v}\|_{L_2([0,1])}^2 = \phi_1^2 a^2 \int_0^1 \cosh^2(\sqrt{\phi_1} x) dx \\ &= \frac{\phi_1^2}{(\sqrt{\phi_1} \sinh(\sqrt{\phi_1}) + \cosh(\sqrt{\phi_1}))^2} \frac{\sqrt{\phi_1} + \cosh(\sqrt{\phi_1}) \sinh(\sqrt{\phi_1})}{2\sqrt{\phi_1}}. \end{aligned}$$

Taking into consideration that

$$\lim_{t \rightarrow \infty} \frac{\cosh(t)}{e^t} = \lim_{t \rightarrow \infty} \frac{\sinh(t)}{e^t} = \frac{1}{2},$$

we obtain for large ϕ_1 :

$$\lim_{\phi_1 \rightarrow \infty} \|\phi_1 \tilde{v}\|_{L_2([0,1])} \phi_1^{-1/4} = \lim_{\phi_1 \rightarrow \infty} \sqrt{\phi_1^{3/2} \frac{e^{\sqrt{\phi_1}}}{2\phi_1 e^{\sqrt{\phi_1}}} \phi_1^{-1/4}} = \frac{1}{\sqrt{2}}.$$

Since $\|f\|_{L_2([0,2])}$ is fixed, we obtain the asymptotics

$$\|\Phi v\|_{L_2([0,2])} \sim \phi_1^{1/4} \|f\|_{L_2([0,2])}.$$

We finally remark that this problem can be lifted by parallel translation to a Poisson problem on $]0, 2[^d$ for $d > 1$, if the newly created boundaries are equipped with homogeneous Neumann boundary conditions. In this case $\partial_{\sigma\nu} v_0(x) = 1$ is constant along the set $\{x \in]0, 2[^d : x_1 = 1\}$, so that our estimate cannot be improved, even if $\partial_{\sigma\nu} v_0$ is assumed to be in a more regular space than $L_2(\partial J_0)$.

Thus, taking also into account our lower bound on the condition number (31) we can be quite sure that the estimates in the following section will be sharp.

5.1 Application to distributed control with control bounds

Let us come back to our elliptic optimal control problem from Section 2.5.1 and consider the case of distributed control, to be solved by a semi-smooth Newton method in function space. We consider preconditioning of the linear systems (4) that arise in each Newton step. In this setting we have $M = E_S : H^1(\Omega) \hookrightarrow L_2(\Omega)$ – the Sobolev embedding, $C^* = \alpha^{-1/2} \chi_{\mathcal{I}}(p) E_S$, and $I = Id$. Thus, we may set $\phi(x) = \alpha^{-1/2} \chi_{\mathcal{I}}(x)$, and thus $\phi_1 = \alpha^{-1/2}$ and $\phi_0 = 0$. Then $(CIMv)(x) = E_S^* \phi(x) (E_S v)(x) = \Phi$ in the notation of (45).

Application of Lemma 5.3 yields the following result:

Proposition 5.5. *Consider the preconditioner $\langle \cdot, \cdot \rangle_Q$ applied to the block operator (4). Assume that the boundary between active and inactive set satisfies the Assumption 5.1. Then we obtain the following condition number:*

$$\kappa_Q \leq c(1 + \alpha^{-1/4}).$$

Proof. In order to apply Lemma 3.6 we set

$$f = M^{-*} (A + CIM)^* v = (A + \Phi)^* v,$$

so that v is the solution of the following problem

$$a(v, w) + \int_{\Omega} \phi(x) v w dx = \int_{\Omega} f w dx \quad \forall w \in H_0^1(\Omega).$$

Hence, Lemma 5.3 yields

$$\|C^*v\|_{L_2(\Omega)} \leq (1 + c\alpha^{-1/8})\|f\|_{L_2(\Omega)}.$$

Inserting this estimate into (28) we obtain the desired result. \square

Remark 5.6. The authors are not aware of any simple *a-priori* assumptions to be imposed on the control constrained problem that would imply Assumption 5.1 for (4) in each Newton step. However, smoothly bounded active sets are observed *a-posteriori* very frequently in numerical examples. Thus, the above theorem explains, why the observed condition numbers of Q are much better than the general estimate ($O(\alpha^{-1/4})$ vs. $O(\alpha^{-1/2})$).

5.2 Application to distributed control with state constraints

For state constraints, we assume that $b(x)$ is piecewise constant taking two values $\|b\|_\infty = b_1 > b_0 > 0$. In penalty methods, as in Example 2.5.4 we have $b(x) = 1 + \gamma\chi_{y < 0}(x)$. As usual in distributed control $U = H = L_2(\Omega)$ and $C = \alpha^{-1/2}E_S^*$, $I = Id$, $M = \sqrt{b(x)}E_S$, so that

$$(CIMv)(x) = E_S^*\alpha^{-1/2}\sqrt{b(x)}(E_Sv)(x).$$

So we can define $\phi(x) := \alpha^{-1/2}\sqrt{b(x)}$ such that $CIM = \Phi$ in the notation of (45). By our assumption, ϕ is also piecewise constant, and we set $\phi_1 = b_1\alpha^{-1/2}$ and $\phi_0 = b_0\alpha^{-1/2}$. Application of Lemma 5.3 then yields:

Proposition 5.7. *Assume that $p \in L_\infty$. Consider the preconditioner $\langle \cdot, \cdot \rangle_Q$ applied to the operator (38). Then we obtain the following condition number:*

$$\kappa_Q \leq c(1 + \|b\|_\infty^{1/4}\alpha^{-1/4}).$$

Proof. We proceed similar as in the control constrained case, defining

$$f = (A + CIM)^*v = (A + \Phi)^*v.$$

Then $\langle v, v \rangle_Q = \|M^{-*}f\|_{L_2(\Omega)}^2$. By Lemma 5.3 we conclude

$$\begin{aligned} \|C^*v\|_{L_2(\Omega)} &= \alpha^{-1/2}\|v\|_{L_2(\Omega)} \leq \alpha^{-1/2}(1 + c\phi_1^{1/4})\|\Phi^{-1}f\|_{L_2(\Omega)} \\ &= \alpha^{-1/2}(1 + c\|b\|_\infty^{1/8}\alpha^{-1/8})\|\sqrt{\alpha}M^{-*}f\|_{L_2(\Omega)} \\ &= (1 + c\|b\|_\infty^{1/8}\alpha^{-1/8})\|M^{-*}f\|_{L_2(\Omega)} \end{aligned}$$

Inserting this into (28) we obtain the desired result. \square

In penalty methods for state constrained problems one considers a homotopy, which results in $\gamma \rightarrow \infty$ and thus $\|b\|_\infty \rightarrow \infty$. In practical applications this leads to values of γ in the order of 10^8 to 10^{12} . Concerning parameters of such high magnitude, our improved condition number estimate is of particular value, since the number of required cg iterations then only grows with the 8th root of $\|b\|_\infty$. As we will see in our numerical examples, compared to the preconditioner Q_0 , which may require hundreds or thousands of cg iterations (proportional to $\sqrt{\|b\|_\infty}$), Q merely takes a very limited number of iterations.

6 Small (in)active sets

Finally, for completeness we will discuss a second case, where the bounds of Section 4 can be improved. In some situations it is likely that those sets where M and C are large have small Lebesgue measure. Applications comprise regularized bang-bang control, where the inactive set becomes small for small regularization parameters, and regularized state constrained control, if the active set tends to a Lebesgue null set.

We approach our problem via an L_∞ -estimate due to Stampacchia:

Lemma 6.1. *Let $\Omega \subset \mathbb{R}^d$ for $d \leq 3$ be a bounded Lipschitz domain and consider the following elliptic equation in weak form:*

$$v \in H_0^1(\Omega) : \quad a(v, w) + \int_{\Omega} \phi(x)vw \, dx = \int_{\Omega} fw \, dx, \quad (55)$$

where $\phi(x)$ is a positive function in $L_\infty(\Omega)$ and $f \in L_2(\Omega)$.

Then we have the estimates

$$\begin{aligned} \|v\|_{L_\infty} &\leq c\|f\|_{L_2} \\ \|v\|_{L_2} &\leq c\|f\|_{L_1}, \end{aligned}$$

where c is independent of ϕ .

Proof. Our first estimate is due to Stampacchia [14, Theorem B.2], and our second estimate follows via a duality technique, similar to the one, used in the proof of Lemma 5.3. \square

Proposition 6.2. *Consider the preconditioners Q_0 from (20) and Q from (24) applied to the block operator (4) derived from the control constrained problem (2). We have the condition number estimates*

$$\begin{aligned} \kappa_{Q_0} &\leq c(1 + \alpha^{-1}\|\chi_{\mathcal{I}}(p)\|_{L_1}) \\ \kappa_Q &\leq c(1 + \alpha^{-1}\|\chi_{\mathcal{I}}(p)\|_{L_1}). \end{aligned}$$

Proof. In both cases we have

$$\langle C^*v, C^*v \rangle = \int_{\Omega} \alpha^{-1}\chi_{\mathcal{I}}(p)v^2 \, dx \leq \|\alpha^{-1}\chi_{\mathcal{I}}(p)\|_{L_1}\|v\|_{L_\infty}^2.$$

By Lemma 6.1, taking into account the definition of A^* , so that $f = M^{-*}A^*v$ according to (55) we obtain

$$\|v\|_{L_\infty}^2 \leq c\|f\|_{L_2}^2 \leq c\langle M^{-*}A^*v, M^{-*}A^*v \rangle_{L_2} = c\langle v, v \rangle_{Q_0}$$

and also $\|v\|_{L_\infty}^2 \leq c\langle v, v \rangle_Q$. Inserting these estimates into (22) and (28), respectively, we obtain the desired result. \square

In bang-bang control, one frequently encounters an assumption of the form (cf. e.g. [5]):

$$|\{x \in \Omega : |p(x)| < \varepsilon\}| \leq c\varepsilon.$$

If such a problem is regularized by a homotopy $\alpha \rightarrow 0$, as in Section 2.5.3, and one assumes that the corresponding adjoint states p_α uniformly satisfy such an assumption, too, then one obtains $\|\chi_{\mathcal{I}}(p_\alpha)\|_{L_1} \leq c\alpha$. Then, the condition number of the preconditioners remains bounded as $\alpha \rightarrow 0$.

For regularized state constraints we have a similar result:

Proposition 6.3. *Consider the preconditioners $\langle \cdot, \cdot \rangle_{Q_0}$ and $\langle \cdot, \cdot \rangle_Q$ applied to the block operator (38) derived from the state constrained problem from Section 2.5.4. We have the condition number estimates*

$$\begin{aligned}\kappa_{Q_0} &\leq c(1 + \alpha^{-1} \|b\|_{L_1}) \\ \kappa_Q &\leq c(1 + \alpha^{-1} \|b\|_{L_1}).\end{aligned}$$

Proof. In both cases we have

$$\langle C^*v, C^*v \rangle \leq \alpha^{-1} \|v\|_{L_2}^2,$$

and it remains to derive a bound of the form

$$\|v\|_{L_2} \leq c \|M^{-*} \underbrace{(A + CIM)^*v}_f\|_{L_2} = c \|M^{-*} f\|_{L_2},$$

for Q and the corresponding one for Q_0 , where the term CIM is missing. In our case we have $(M^*w)(x) = \sqrt{b(x)}w(x)$. Since $(A + CIM)^*v = f$, we conclude with Lemma 6.1 that

$$\|v\|_{L_2} \leq \|f\|_{L_1} = \|\sqrt{b}M^{-*}f\|_{L_1} \leq \|\sqrt{b}\|_{L_2} \|M^{-*}f\|_{L_2} \leq \sqrt{\|b\|_{L_1}} \|M^{-*}f\|_{L_2}.$$

Inserting this estimate (and the corresponding one for Q_0) into (22) and (28), respectively, we obtain the desired result. \square

7 Numerical examples

In this section we perform a numerical study of the pcg method with our preconditioners. We consider two examples. The first is related to a control constrained problem, the second is related to a regularized state constrained problem.

The pcg iteration is terminated, after the estimated error in energy norm has dropped below 10^{-8} , where an estimator in the spirit of [6, Sec. 5.3.3(c)] is used.

Our simple implementation is based on matlab, and the discretization of our optimality system was done with 5-point star finite differences. For the solution of the single PDE blocks, the built-in sparse direct Cholesky factorization has been used. In very large scale applications, in particular in the 3d case, the sparse direct solver will have to be replaced by a properly preconditioned iterative solver, of course.

In both cases the computational domain is the unit-square, and the current active set \mathcal{A} is a disc with center $(0.5, 0.5)$ and radius 0.4. As right hand side we choose the function $r \equiv 1$.

In our first problem, we thus solve a problem of the form

$$\begin{pmatrix} E^*E & A^* \\ A & -E^*\alpha^{-1}\chi_{\Omega \setminus \mathcal{A}}E \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad (56)$$

with varying Tychonov parameter α . Here A corresponds to the weak form of $-\Delta$ on $H_0^1([0, 1] \times [0, 1])$, and E is the Sobolev embedding. Of particular interest is the case, where α is very small. In the second example we solve a problem of the form

$$\begin{pmatrix} E^*(1 + \gamma\chi_{\mathcal{A}})E & A^* \\ A & -E^*E \end{pmatrix} \begin{pmatrix} \delta y \\ \delta p \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad (57)$$

$h \setminus \alpha$	10^{-2}	10^{-4}	10^{-6}	10^{-8}	10^{-10}	10^{-12}
2^{-6}	8	13	18	22	24	23
2^{-7}	7	13	18	22	29	32
2^{-8}	7	13	17	22	31	41
2^{-9}	7	12	17	22	31	52
2^{-10}	7	12	17	22	33	54

$h \setminus \alpha$	10^{-2}	10^{-4}	10^{-6}	10^{-8}	10^{-10}	10^{-12}
2^{-6}	5	9	33	234	798	953
2^{-7}	5	9	33	236	1749	3222
2^{-8}	5	9	31	233	2095	> 10000
2^{-9}	4	9	31	231	2122	> 10000

Figure 1: Number of pcg iterations for preconditioner Q (top) and Q_0 (bottom) for control constraints with varying grid size h and Tychonov parameter α .

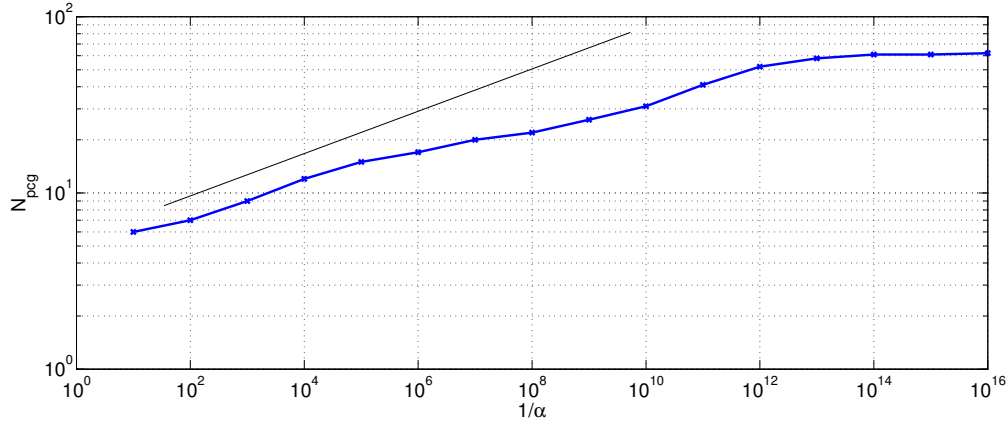


Figure 2: Number N_{pcg} of computed pcg iterations for problem (56) and theoretical prediction $N_{pcg} \sim \alpha^{-1/8}$ plotted against α^{-1} on a 512×512 grid ($h = 2^{-9}$).

with varying penalty parameter γ . Here γ can become very large during the course of a path-following method. It can be observed in both examples that the preconditioner Q is vastly superior to Q_0 for large α^{-1} or γ . Moreover, for very large parameters discretization effects tend to yield smaller numbers of pcg iterations for coarse grids than for fine grids.

If we compare the observed number of pcg iterations with the predicted number of iterations (cf. Figure 7) we observe two things. First, the average increase of iterations seems to be slightly slower than predicted. However, there are regions (i.e. $\alpha^{-1} \in [10^2, 10^4]$ and $\alpha^{-1} \in [10^{10}, 10^{12}]$) where the slopes seem to fit. For very large $\alpha^{-1} > 10^{12}$ we observe a saturation of the number of iterations. This effect is most probably due to the discretization of the problem, and can also be observed in the last columns of the top table in Figure 1. Such behavior is usually a clear indication that the discretization of the problem is too coarse.

Of course, this iterative solver can also be used as an inner loop inside a semi-smooth Newton method. Let us give an example in the context of control constraints $u \geq 0$. Here we use again a control problem with distributed control and corresponds to the weak form of

$h \setminus \gamma$	10^2	10^4	10^6	10^8	10^{10}	10^{12}
2^{-6}	8	12	17	23	28	29
2^{-7}	8	12	17	23	34	38
2^{-8}	8	12	16	23	36	49
2^{-9}	8	12	16	23	37	58
2^{-10}	7	12	16	23	36	61

$h \setminus \gamma$	10^2	10^4	10^6	10^8	10^{10}	10^{12}
2^{-6}	5	11	41	300	1204	1638
2^{-7}	5	11	39	297	2464	5323
2^{-8}	5	11	39	293	2727	> 10000
2^{-9}	5	11	40	293	2722	> 10000

Figure 3: Number of pcg iterations for preconditioner Q (top) and Q_0 (bottom) for regularized state constraints with varying grid size h and penalty parameter γ .

$\alpha \setminus$	N_{Newton}	$N_{\text{pcg}}^{\text{total}}$	$N_{\text{pcg}}^{\text{avg}}$
10^{-2}	3	8	2.667
10^{-4}	6	28	4.667
10^{-6}	10	59	5.9
10^{-8}	9	79	8.778
10^{-10}	9	119	13.22
10^{-12}	8	142	17.75

Figure 4: Convergence history of semi-smooth Newton method for varying α

$\mathcal{A} = -\Delta$ on $H_0^1([0, 1] \times [0, 1])$. As desired state, we choose $y_d = A^{-1} \sin(\pi x_1 x_2)$. State y and adjoint state p are both discretized by standard finite differences with mesh size $h = 2^{-8}$. To compute the solution for this problem with varying, up to very small α we reuse the computed solution for the last (larger) α as an initial guess for the next (smaller) α . This acts as a simplistic path-following method for $\alpha \rightarrow 0$ and compensates for the inefficient global convergence behavior of semi-smooth Newton in case of small α .

8 Conclusion and Outlook

We have proposed and analyzed block preconditioners for systems that arise in certain optimal control problems with PDEs. It can be used effectively for control constraints, if the domain of observation contains the domain of control. For state constraints the approximate active constraint set should be contained in the control domain.

In these cases the condition numbers of the resulting systems are in general only the square root of the condition numbers, that are obtained via a simple preconditioner Q_0 , asymptotically with respect to critical parameters. Under additional structural assumptions, it can be shown that the condition number grows even slower, like the fourth root of the condition number of Q_0 . In the unconstrained case one obtains condition numbers independent of the critical parameter.

Our results are of particular interest in state constrained optimal control problems, where up to now the robustness of available preconditioners with respect to regularization

parameters was poor. The class of state constrained problems seems to be divided into two subclasses. The first, where the control can act on the active set of the constraints, and the class of remaining problems, where the control acts more indirectly. The first class seems to be tractable more easily than the second, and new ideas are needed for the second class.

However, much work remains to be done. Up to now, our results are only valid for exact solutions of the modified PDEs. This is already a significant progress, since direct solvers for elliptic problems are much more efficient than for coupled systems. In the parabolic case the advantage is even larger, because the differential equations can be solved by time-stepping procedures. This is much easier than solving the complete coupled system. It is a straightforward idea to replace direct solvers by multigrid preconditioners. However, showing optimality and robustness of these preconditioners seems an open non-trivial theoretical issue that needs to be addressed in the future. The usual H^1 techniques cannot be used because the natural space for the preconditioners is D_K . From an algorithmic point of view it is desirable to apply our preconditioners in an adaptive multilevel method in the spirit of [21], where inexactness of Newton steps caused by the iterative solver and the adaptive grid refinement is appropriately handled within an adaptive inexact Newton path-following method in function space. Finally, it remains to investigate, how much of our results can be applied to more general classes of PDEs, say from continuum mechanics or reaction-diffusion systems.

Acknowledgement The authors would like to thank the referees, who with their detailed remarks contributed significantly to the quality of the presentation.

References

- [1] A. Battermann and M. Heinkenschloss. Preconditioners for Karush-Kuhn-Tucker matrices arising in the optimal control of distributed systems. In *Control and estimation of distributed parameter systems (Vorau, 1996)*, volume 126 of *Internat. Ser. Numer. Math.*, pages 15–32. Birkhäuser, 1998.
- [2] A. Battermann and E. W. Sachs. Block preconditioners for KKT systems in PDE-governed optimal control problems. In *Fast solution of discretized optimization problems (Berlin, 2000)*, volume 138 of *Internat. Ser. Numer. Math.*, pages 1–18. Birkhäuser, Basel, 2001.
- [3] A. Borzi. Smoothers for control- and state-constrained optimal control problems. *Comput. Vis. Sci.*, 11(1):59–66, 2008.
- [4] A. Borzi and V. Schulz. Multigrid methods for PDE optimization. *SIAM Rev.*, 51(2):361–395, 2009.
- [5] K. Deckelnick and M. Hinze. A note on the approximation of elliptic control problems with bang-bang controls. *Comput. Optim. Appl.*, pages 931–939, 2012.
- [6] P. Deuffhard and M. Weiser. *Adaptive Numerical Solution of Partial Differential Equations*. de Gruyter, 2012.
- [7] H. S. Dollar, N. I. M. Gould, M. Stoll, and A. J. Wathen. Preconditioning saddle-point systems with applications in optimization. *SIAM J. Sci. Comput.*, 32(1):249–270, 2010.

- [8] H. Gajewski, K. Gröger, and K. Zacharias. *Nichtlineare Operatorgleichungen und Operatordifferentialgleichungen*. Akademie-Verlag, Berlin, 1974. Mathematische Lehrbücher und Monographien, II. Abteilung, Mathematische Monographien, Band 38.
- [9] A. Günnel, R. Herzog, and E. Sachs. A note on preconditioners and scalar products for Krylov methods in Hilbert space. Technical report, TU Chemnitz, July 2011.
- [10] R. Herzog and E. Sachs. Preconditioned conjugate gradient method for optimal control problems with control and state constraints. *SIAM J. Matrix Anal. Appl.*, 31(5):2291–2317, 2010.
- [11] M. Hintermüller, K. Ito, and K. Kunisch. The primal-dual active set strategy as a semi-smooth Newton method. *SIAM J. Optim.*, 13:865–888, 2003.
- [12] M. Hintermüller and K. Kunisch. Feasible and non-interior path-following in constrained minimization with low multiplier regularity. *SIAM J. Control Optim.*, 45(4):1198–1221, 2006.
- [13] M. Hinze. A variational discretization concept in control constrained optimization: the linear-quadratic case. *Comput. Optim. Appl.*, 30:45–63, 2005.
- [14] D. Kinderlehrer and G. Stampacchia. *An Introduction to Variational Inequalities and their Applications*. Academic Press, New York, 1980.
- [15] K. Krumbiegel and A. Rösch. A virtual control concept for state constrained optimal control problems. *Comput. Optim. Appl.*, 43(2):213–233, 2009.
- [16] J. L. Lions and E. Magenes. *Nonhomogeneous Boundary Value Problems and Applications*, volume 1–3. Springer-Verlag, Berlin, 1972.
- [17] J. W. Pearson, M. Stoll, and A. J. Wathen. Preconditioners for state constrained optimal control problems with Moreau-Yosida penalty function. Technical Report NA-12/05, Oxford University, Mathematical Institute, March.
- [18] T. Rees, M. Stoll, and A. Wathen. All-at-once preconditioning in PDE-constrained optimization. *Kybernetika (Prague)*, 46(2):341–360, 2010.
- [19] A. Schiela. A simplified approach to semismooth Newton methods in function space. *SIAM J. Optim.*, 19(3):1417–1432, 2008.
- [20] A. Schiela. Barrier methods for optimal control problems with state constraints. *SIAM J. Optim.*, 20(2):1002–1031, 2009.
- [21] A. Schiela and A. Günther. An interior point algorithm with inexact step computation in function space for state constrained optimal control. *Numer. Math.*, 119(2):373–407, 2011.
- [22] J. Schöberl, R. Simon, and W. Zulehner. A robust multigrid method for elliptic optimal control problems. *SIAM J. Numer. Anal.*, 49(4):1482–1503, 2011.
- [23] J. Schöberl and W. Zulehner. Symmetric indefinite preconditioners for saddle point problems with applications to PDE-constrained optimization problems. *SIAM J. Matrix Anal. Appl.*, 29(3):752–773, 2007.

- [24] M. Stoll and A. Wathen. Preconditioning for partial differential equation constrained optimization with control constraints. *Numerical Linear Algebra with Applications*, 19(1):53–71, 2012.
- [25] F. Tröltzsch. *Optimal Control of Partial Differential Equations: Theory, Methods and Applications*. AMS, Providence, 2010.
- [26] M. Ulbrich. Semismooth Newton methods for operator equations in function spaces. *SIAM J. Optim.*, 13:805–842, 2003.
- [27] D. Wachsmuth and G. Wachsmuth. Necessary conditions for convergence rates of regularizations of optimal control problems. RICAM Report 2012-04, Johann Radon Institute for Computational and Applied Mathematics, January 2012.
- [28] W. Zulehner. Nonstandard norms and robust estimates for saddle point problems. *SIAM J. Matrix Anal. Appl.*, 32(2):536–560, 2011.