

## Part XI, Chapter 50

---

### Mixed finite element approximation

This chapter is concerned with the approximation of the model problem analyzed in Chapter 49. We focus on the Galerkin approximation in the conforming setting. We establish necessary and sufficient conditions for well-posedness, and we derive error bounds in terms of the best-approximation error. Then we consider the algebraic viewpoint, and we discuss augmented Lagrangian methods in the context of saddle point problems. Finally, we examine iterative solvers, including Uzawa iterations and Krylov subspace methods.

#### 50.1 Conforming Galerkin approximation

Let  $V$  and  $Q$  be two reflexive (real) Banach spaces. Let  $a$  and  $b$  be two bounded bilinear forms on  $V \times V$  and on  $V \times Q$  respectively. Let  $f \in V'$  and let  $g \in Q'$ . We consider the following model problem:

$$\begin{cases} \text{Find } u \in V \text{ and } p \in Q \text{ such that} \\ a(u, w) + b(w, p) = f(w), & \forall w \in V, \\ b(u, q) = g(q), & \forall q \in Q. \end{cases} \quad (50.1)$$

We introduce the associated operators  $A \in \mathcal{L}(V; V')$  and  $B \in \mathcal{L}(V; Q')$  such that  $a(v, w) := \langle A(v), w \rangle_{V', V}$  and  $b(v, q) := \langle B(v), q \rangle_{Q', Q}$ . We assume that (50.1) is well-posed. Owing to Theorem 49.13, this means that the bilinear form  $a$  satisfies the conditions (49.36) (implying the inf-sup condition  $\inf_{v \in \ker(B)} \sup_{w \in \ker(B)} \frac{|a(v, w)|}{\|v\|_V \|w\|_V} =: \alpha > 0$ ) and that the bilinear form  $b$  satisfies the inf-sup condition (49.37), i.e.,  $\inf_{q \in Q} \sup_{v \in V} \frac{|b(v, q)|}{\|v\|_V \|q\|_Q} =: \beta > 0$ .

A conforming Galerkin approximation of (50.1) is obtained by considering two finite-dimensional subspaces  $V_h \subset V$ ,  $Q_h \subset Q$ . The discrete problem is

$$\begin{cases} \text{Find } u_h \in V_h \text{ and } p_h \in Q_h \text{ such that} \\ a(u_h, w_h) + b(w_h, p_h) = f(w_h), & \forall w_h \in V_h, \\ b(u_h, q_h) = g(q_h), & \forall q_h \in Q_h. \end{cases} \quad (50.2)$$

### 50.1.1 Well-posedness

Let  $B_h : V_h \rightarrow Q'_h$  be the discrete counterpart of the operator  $B : V \rightarrow Q'$ , that is,  $\langle B_h(v_h), q_h \rangle_{Q'_h, Q_h} := \langle B(v_h), q_h \rangle_{Q', Q} = b(v_h, q_h)$  for all  $(v_h, q_h) \in V_h \times Q_h$ . The null space of  $B_h$  is such that

$$\ker(B_h) = \{v_h \in V_h \mid \forall q_h \in Q_h, b(v_h, q_h) = 0\}. \quad (50.3)$$

One important aspect of the discretization is that the surjectivity of  $B$  does not imply that of  $B_h$ . One rare occasion where this is nevertheless the case is when  $B^*(Q_h) \subset V_h$ , i.e.,  $B_h^* = B^*|_{Q_h}$ . This exceptional situation is illustrated in Exercise 49.6. Note also that in general,  $\ker(B_h)$  is not necessarily a subspace of  $\ker(B)$ .

**Proposition 50.1 (Well-posedness).** (50.2) is well-posed if and only if

$$\inf_{v_h \in \ker(B_h)} \sup_{w_h \in \ker(B_h)} \frac{|a(v_h, w_h)|}{\|v_h\|_V \|w_h\|_V} := \alpha_h > 0, \quad (50.4a)$$

$$\inf_{q_h \in Q_h} \sup_{v_h \in V_h} \frac{|b(v_h, q_h)|}{\|v_h\|_V \|q_h\|_Q} := \beta_h > 0. \quad (50.4b)$$

*Proof.* Apply Theorem 49.13 and use the fact that (50.4a) implies both conditions in (49.36) since  $V_h$  is finite-dimensional; see Remark 26.7.  $\square$

The condition (50.4a) holds true for all conforming subspaces  $V_h$  and  $Q_h$  if  $a$  is  $V$ -coercive (the coercivity of  $a$  on  $\ker(B)$  may not be sufficient since it may happen that  $\ker(B_h) \not\subset \ker(B)$ ). Note that verifying the inf-sup condition for  $a$  on  $V_h \times V_h$  is not sufficient to prove (50.4a) (think of an invertible matrix having a square diagonal sub-block that is not invertible; see Exercise 50.1). Furthermore, the condition (50.4b) is equivalent to  $B_h$  being surjective, which is also equivalent to  $B_h^*$  being injective since the setting is finite-dimensional. In practice, it is important that both (50.4a) and (50.4b) hold true uniformly w.r.t.  $h \in \mathcal{H}$ , i.e.,  $\inf_{h \in \mathcal{H}} \alpha_h =: \alpha_0 > 0$  and  $\inf_{h \in \mathcal{H}} \beta_h =: \beta_0 > 0$ .

### 50.1.2 Error analysis

Our goal is to estimate the approximation errors  $(u - u_h)$  and  $(p - p_h)$  in terms of the best-approximation error on  $u$  by a discrete field in  $V_h$  and the best-approximation error on  $p$  by a discrete function in  $Q_h$ . Céa's lemma (Lemma 26.13) could be applied to the bilinear form  $t((v, q), (w, r)) := a(v, w) + b(w, q) + b(v, r)$  introduced in §49.4.2 (see Exercise 50.2). But here, we present a more specific analysis distinguishing the errors on  $u$  and on  $p$ .

We say that  $\Pi_h \in \mathcal{L}(V; V_h)$  is a Fortin operator for the bilinear form  $b$  if  $b(\Pi_h(v) - v, q_h) = 0$  for all  $q_h \in Q_h$  and all  $v \in V$ . (We do not assume  $V_h$  to be pointwise invariant under  $\Pi_h$ .) This class of operators is investigated in §26.2.3. In particular, Lemma 26.9 shows that the inf-sup condition (50.4b) implies the existence of a Fortin operator with  $\|\Pi_h\|_{\mathcal{L}(V; V_h)} \leq \frac{\|b\|}{\beta_h}$ .

**Lemma 50.2 (Error estimate).** *Let  $(u, p)$  solve (50.1). Assume (50.4) and let  $(u_h, p_h)$  solve the discrete problem (50.2). Let  $\Pi_h \in \mathcal{L}(V; V_h)$  be any Fortin operator. The following error estimates hold true:*

$$\|u - u_h\|_V \leq c_{1h} \inf_{v_h \in \Pi_h(u) + \ker(B_h)} \|u - v_h\|_V + c_{2h} \inf_{q_h \in Q_h} \|p - q_h\|_Q, \quad (50.5a)$$

$$\|p - p_h\|_Q \leq c_{3h} \inf_{v_h \in \Pi_h(u) + \ker(B_h)} \|u - v_h\|_V + c_{4h} \inf_{q_h \in Q_h} \|p - q_h\|_Q, \quad (50.5b)$$

with  $c_{1h} := (1 + \frac{\|a\|}{\alpha_h})$ ,  $c_{2h} := \frac{\|b\|}{\alpha_h}$  if  $\ker(B_h) \not\subset \ker(B)$  and  $c_{2h} := 0$  otherwise,  $c_{3h} := c_{1h} \frac{\|a\|}{\beta_h}$ , and  $c_{4h} := 1 + \frac{\|b\|}{\beta_h} + c_{2h} \frac{\|a\|}{\beta_h}$ .

*Proof.* (1) Estimate on  $(u - u_h)$ . Let  $v_h \in \Pi_h(u) + \ker(B_h)$ , i.e.,  $v_h := \Pi_h(u) + \gamma_h$  with  $\gamma_h \in \ker(B_h)$ . Then  $u_h - v_h \in \ker(B_h)$  since we have

$$b(u_h - v_h, q_h) = b(u_h - \Pi_h(u), q_h) + b(\gamma_h, q_h) = b(u_h - u, q_h) = 0,$$

for all  $q_h \in Q_h$ , where we used the Galerkin orthogonality property for the second equation in (50.2). Owing the inf-sup condition (50.4a), we infer that

$$\begin{aligned} \alpha_h \|u_h - v_h\|_V &\leq \sup_{y_h \in \ker(B_h)} \frac{|a(u_h - v_h, y_h)|}{\|y_h\|_V} \\ &= \sup_{y_h \in \ker(B_h)} \frac{|b(y_h, p - p_h) + a(u - v_h, y_h)|}{\|y_h\|_V}, \end{aligned}$$

where the equality follows from the Galerkin orthogonality property for the first equation in (50.2), i.e., we have  $a(u - u_h, y_h) + b(y_h, p - p_h) = 0$  for all  $y_h \in V_h$ . If  $\ker(B_h) \subset \ker(B)$ , then  $b(y_h, p - p_h) = 0$  for all  $y_h \in \ker(B_h)$ , yielding

$$\alpha_h \|u_h - v_h\|_V \leq \|a\| \|u - v_h\|_V.$$

In the general case, we have  $b(y_h, p_h) = 0 = b(y_h, q_h)$  for all  $q_h \in Q_h$ , since  $y_h$  is in  $\ker(B_h)$ . This implies that

$$\alpha_h \|u_h - v_h\|_V \leq \|a\| \|u - v_h\|_V + \|b\| \|p - q_h\|_Q.$$

Hence, both cases are summarized by the following estimate:

$$\|u_h - v_h\|_V \leq \frac{\|a\|}{\alpha_h} \|u - v_h\|_V + c_{2h} \|p - q_h\|_Q,$$

with  $c_{2h}$  as in the assertion. Using the triangle inequality and taking the infimum over  $v_h \in \Pi_h(u) + \ker(B_h)$  and over  $q_h \in Q_h$  leads to (50.5a).

(2) Estimate on  $(p - p_h)$ . Using again the Galerkin orthogonality property for the first equation in (50.2), we have

$$b(v_h, q_h - p_h) = a(u_h - u, v_h) + b(v_h, q_h - p), \quad \forall (v_h, q_h) \in V_h \times Q_h.$$

Combined with the inf-sup condition (50.4b), this implies that

$$\begin{aligned} \beta_h \|q_h - p_h\|_Q &\leq \sup_{v_h \in V_h} \frac{|b(v_h, q_h - p_h)|}{\|v_h\|_V} \\ &\leq \|a\| \|u - u_h\|_V + \|b\| \|p - q_h\|_Q. \end{aligned}$$

The bound (50.5b) follows from the triangle inequality, the bound on  $(u - u_h)$ , and by taking the infimum over  $q_h \in Q_h$ .  $\square$

The estimate on  $(u - u_h)$  involves the best-approximation error on  $u$  by a member of the affine subspace  $\Pi_h(u) + \ker(B_h)$ . This error may not be easy to estimate in practice, and it is sometimes preferable to bound it by the best-approximation error on  $u$  by a member of  $V_h$  since  $\Pi_h(u) + \ker(B_h) \subset V_h$ . Of course, the best-approximation error in  $\Pi_h(u) + \ker(B_h)$  is larger than the best-approximation error in  $V_h$ . The following lemma quantifies the discrepancy. (Recall that (50.4b) is equivalent to the existence of Fortin operators.)

**Lemma 50.3 (Best-approximation in  $V_h$ ).** *Assume (50.4b). The following holds true for all  $u \in V$  and any Fortin operator  $\Pi_h \in \mathcal{L}(V; V_h)$ :*

$$\inf_{v_h \in \Pi_h(u) + \ker(B_h)} \|u - v_h\|_V \leq (1 + \|\Pi_h\|_{\mathcal{L}(V; V_h)}) \inf_{y_h \in V_h} \|u - y_h\|_V. \quad (50.6)$$

*Proof.* Let  $u \in V$ . Let  $y_h \in V_h$  and set  $z_h := \Pi_h(u - y_h)$ . Then  $y_h + z_h = \Pi_h(u) + y_h - \Pi_h(y_h) \in \Pi_h(u) + \ker(B_h)$  since  $b(y_h - \Pi_h(y_h), q_h) = 0$  for all  $q_h \in Q_h$ . This implies that

$$\begin{aligned} \inf_{v_h \in \Pi_h(u) + \ker(B_h)} \|u - v_h\|_V &\leq \|u - (y_h + z_h)\|_V \leq \|u - y_h\|_V + \|z_h\|_V \\ &\leq (1 + \|\Pi_h\|_{\mathcal{L}(V; V_h)}) \|u - y_h\|_V, \end{aligned}$$

and we conclude by taking the infimum over  $y_h \in V_h$ .  $\square$

**Remark 50.4 ( $g = 0$ ).** If  $g = 0$ , then  $\Pi_h(u) \in \ker(B_h)$ , and the infimum in (50.5) and (50.6) reduces to  $v_h \in \ker(B_h)$ .  $\square$

**Corollary 50.5 (Error estimate).** *Let  $(u, p)$  solve (50.1). Assume (50.4) and let  $(u_h, p_h)$  solve the discrete problem (50.2). The following error estimates hold true:*

$$\|u - u_h\|_V \leq c'_{1h} \inf_{v_h \in V_h} \|u - v_h\|_V + c_{2h} \inf_{q_h \in Q_h} \|p - q_h\|_Q, \quad (50.7a)$$

$$\|p - p_h\|_Q \leq c'_{3h} \inf_{v_h \in V_h} \|u - v_h\|_V + c_{4h} \inf_{q_h \in Q_h} \|p - q_h\|_Q, \quad (50.7b)$$

with  $c'_{1h} := (1 + \frac{\|a\|}{\alpha_h})(1 + \|\Pi_h\|_{\mathcal{L}(V;V_h)})$  for every Fortin operator  $\Pi_h \in \mathcal{L}(V;V_h)$ ,  $c'_{3h} := c'_{1h} \frac{\|a\|}{\beta_h}$ , and  $c_{2h}, c_{4h}$  are as in Lemma 50.2.

*Proof.* Combine Lemma 50.2 with Lemma 50.3.  $\square$

**Remark 50.6** ( $c'_{1h}$ ). Lemma 26.9 asserts the existence of a Fortin operator with  $\|\Pi_h\|_{\mathcal{L}(V;V_h)} \leq \frac{\|b\|}{\beta_h}$ . Hence, the upper bound  $c'_{1h} \leq (1 + \frac{\|a\|}{\alpha_h})(1 + \frac{\|b\|}{\beta_h})$  always holds true. However, the estimate  $\|\Pi_h\|_{\mathcal{L}(V;V_h)} \leq \frac{\|b\|}{\beta_h}$  can be pessimistic. For instance, for the Stokes equations in elongated domains, the boundedness constant of the bilinear form  $b(\mathbf{v}, p) = (\nabla \cdot \mathbf{v}, p)_{L^2(D)}$  on  $\mathbf{H}_0^1(D) \times L^2(D)$  is  $\|b\| = 1$ , and the inf-sup constant  $\beta_h$  can be shown to be very small (see Chizhonkov and Olshanskii [119], Dobrowolski [170]), whereas for some of these domains it is possible to construct a Fortin operator with norm of order unity (see Mardal et al. [294], Linke et al. [284]).  $\square$

**Remark 50.7** ( $\ker(B_h)$ ). We refer the reader to Theorem 51.16 for an example of error estimate exploiting the approximation properties in  $\ker(B_h)$  in the context of Darcy's equations.  $\square$

**Remark 50.8** ( $c_{2h}$ ). The constant  $c_{2h}$  vanishes whenever  $\ker(B_h) \subset \ker(B)$ . Using a discrete pair  $(V_h, Q_h)$  that guarantees that  $\ker(B_h) \subset \ker(B)$  may be interesting when the best approximation error on  $p$  is (much) larger than that on  $u$ . A simple example where this occurs is when  $f = B^*(p)$  for some  $p \in Q$  and  $g = 0$ , so that the solution to (50.1) is  $(0, p)$ . If  $\ker(B_h) \subset \ker(B)$ , the estimate (50.7a) implies that  $u_h = u = 0$ . But if  $\ker(B_h) \not\subset \ker(B)$ , then  $u_h$  is generally nonzero and grows linearly with the size of  $p$ , which is not a desirable property. More generally, the well-posedness of (50.1) with  $g := 0$  implies the abstract *Helmholtz decomposition*  $V' = Y_0 \oplus Y_1$  with  $Y_0 := A(\ker(B))$  and  $Y_1 = \text{im}(B^*)$ . Whenever the component of  $f$  in  $Y_1$  is much larger than that in  $Y_0$ , the best-approximation error on  $p$  dominates the approximation error on  $u$  unless the discretization satisfies  $\ker(B_h) \subset \ker(B)$ . See also Remark 53.22 for further insight in the context of the Stokes equations.  $\square$

**Remark 50.9 (Stabilization)**. It is possible to approximate (50.1) using discrete spaces  $V_h$  and  $Q_h$  that violate the inf-sup condition (50.4b) by replacing the bilinear forms  $a$  and  $b$  by some stabilized versions  $a_h$  and  $b_h$ ; see Chapters 62 and 63.  $\square$

We now establish an error estimate on  $u$  in a norm that is weaker than that in  $V$ . We do so by using a duality argument in the spirit of the Aubin–Nitsche lemma (Lemma 32.11).

**Definition 50.10 (Smoothing property).** *The problem (50.1) is said to have a smoothing property if there is a Hilbert space  $H \hookrightarrow V$  with inner product  $(\cdot, \cdot)_H$ , two Banach spaces  $Y \hookrightarrow V$  and  $N \hookrightarrow Q$ , and a constant  $c_{\text{smo}}$  such that the following adjoint problem:*

$$\begin{cases} \text{Find } \varphi(g) \in V \text{ and } \vartheta(g) \in Q \text{ such that} \\ a(v, \varphi(g)) + b(v, \vartheta(g)) = (g, v)_H, & \forall v \in V, \\ b(\varphi(g), q) = 0, & \forall q \in Q, \end{cases}$$

has a unique solution for all  $g \in H$  and satisfies the a priori estimate  $\|\varphi(g)\|_Y + \|\vartheta(g)\|_N \leq c_{\text{smo}} \|g\|_H$ .

In addition to the smoothing property, we assume that the spaces  $H$ ,  $Y$ , and  $N$  satisfy an additional approximation property, i.e., there are  $s > 0$  and  $c$  such that the following holds true for all  $(v, q) \in Y \times N$  and all  $h \in \mathcal{H}$ :

$$\inf_{(v_h, q_h) \in V_h \times Q_h} (\|v - v_h\|_V + \|q - q_h\|_Q) \leq c h^s (\|v\|_Y + \|q\|_N). \quad (50.8)$$

**Lemma 50.11 (Improved error estimate in weaker norm).** *Let  $(u, p)$  solve (50.1). Assume (50.4) and let  $(u_h, p_h)$  solve the discrete problem (50.2). Assume that (50.1) has a smoothing property and that (50.8) holds true. Then we have*

$$\|u - u_h\|_H \leq c h^s (\|u - u_h\|_V + \|p - p_h\|_Q),$$

where  $c$  is independent of  $(u, p)$ ,  $(u_h, p_h)$  and  $h \in \mathcal{H}$ .

*Proof.* Set  $V := V \times Q$ ,  $Z := Y \times N$ , and  $L := H \times Q$ , each equipped with the product norm. Define the symmetric positive bilinear form  $l((v, q), (w, r)) := (v, w)_H$  and the seminorm  $|(v, q)|_L := \|v\|_H$ . Apply Lemma 32.11 in the conforming setting with the bilinear form  $t((u, p), (v, q)) := a(u, v) + b(v, p) + b(u, q)$  to conclude.  $\square$

## 50.2 Algebraic viewpoint

In this section, we study the linear system associated with the discrete problem (50.2) assuming that the well-posedness conditions (50.4a)-(50.4b) are satisfied. We also assume that the bilinear form  $a$  satisfies an inf-sup condition on  $V_h \times V_h$ . For simplicity, we consider real vector spaces.

### 50.2.1 The coupled linear system

Let  $N := \dim(V_h)$  and  $M := \dim(Q_h)$ . Let  $\{\varphi_i\}_{i \in \{1:N\}}$  be a basis for  $V_h$  and let  $\{\psi_k\}_{k \in \{1:M\}}$  be a basis for  $Q_h$ . Recall that these bases consist of global shape functions when  $V_h$  and  $Q_h$  are finite element spaces. Proceeding as in §28.1.1, for every column vectors  $\mathbf{U} = (\mathbf{U}_1, \dots, \mathbf{U}_N)^T$  in  $\mathbb{R}^N$

and  $\mathbf{P} = (\mathbf{P}_1, \dots, \mathbf{P}_M)^\top$  in  $\mathbb{R}^M$ , we define the functions  $\mathbf{R}_\varphi(\mathbf{U}) \in V_h$  and  $\mathbf{R}_\psi(\mathbf{P}) \in Q_h$  by  $\mathbf{R}_\varphi(\mathbf{U}) := \sum_{i \in \{1:N\}} \mathbf{U}_i \varphi_i$  and  $\mathbf{R}_\psi(\mathbf{P}) := \sum_{k \in \{1:M\}} \mathbf{P}_k \psi_k$ . The correspondences between  $\mathbf{R}_\varphi(\mathbf{U})$  and  $\mathbf{U}$  and between  $\mathbf{R}_\psi(\mathbf{P})$  and  $\mathbf{P}$  are one-to-one since  $\{\varphi_i\}_{i \in \{1:N\}}$  and  $\{\psi_k\}_{k \in \{1:M\}}$  are bases.

Inserting the expansions of  $\mathbf{R}_\varphi(\mathbf{U})$  and  $\mathbf{R}_\psi(\mathbf{P})$  into (50.2) and choosing the basis functions of  $V_h$  and  $Q_h$  to test (50.2), we obtain the linear system

$$\mathcal{C} \begin{pmatrix} \mathbf{U} \\ \mathbf{P} \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{G} \end{pmatrix}, \quad \mathcal{C} := \begin{pmatrix} \mathcal{A} & \mathcal{B}^\top \\ \mathcal{B} & \mathbf{0} \end{pmatrix}, \quad (50.9)$$

where the matrices  $\mathcal{A} \in \mathbb{R}^{N \times N}$  and  $\mathcal{B} \in \mathbb{R}^{M \times N}$  are such that  $\mathcal{A}_{ij} := a(\varphi_j, \varphi_i)$  and  $\mathcal{B}_{ki} := b(\varphi_i, \psi_k)$  for all  $i \in \{1:N\}$  and all  $k \in \{1:M\}$ ,  $\mathbf{0}$  is the zero matrix in  $\mathbb{R}^{M \times M}$ , and the vectors  $\mathbf{F} \in \mathbb{R}^N$  and  $\mathbf{G} \in \mathbb{R}^M$  are such that  $\mathbf{F}_i = f(\varphi_i)$  and  $\mathbf{G}_k = g(\psi_k)$  for all  $i, j \in \{1:N\}$  and all  $k \in \{1:M\}$ .

The matrix  $\mathcal{C}$  is invertible since (50.2) is well-posed owing to (50.4a)-(50.4b). Notice also that (50.4b) implies that  $\mathcal{B}^\top$  has full column rank and  $\mathcal{B}$  has full row rank (these ranks are equal to  $M$ ). Moreover,  $\mathcal{A}$  is invertible since we additionally assumed that  $a$  satisfies an inf-sup condition on  $V_h \times V_h$ . Algebraic counterparts of the boundedness and inf-sup conditions on the bilinear forms  $a$  and  $b$  can be established by using the dual norm

$$\|\mathbf{U}\|_{\ell_\varphi^2} := \sup_{\mathbf{Y} \in \mathbb{R}^N} \frac{\mathbf{U}^\top \mathbf{Y}}{\|\mathbf{R}_\varphi(\mathbf{Y})\|_V}, \quad \forall \mathbf{U} \in \mathbb{R}^N. \quad (50.10)$$

**Proposition 50.12 (Norm equivalence).** *The following holds true:*

$$\alpha_h \|\mathbf{R}_\varphi(\mathbf{U})\|_V \leq \|\mathcal{A}\mathbf{U}\|_{\ell_\varphi^2} \leq \|a\| \|\mathbf{R}_\varphi(\mathbf{U})\|_V, \quad \forall \mathbf{U} \in \mathbb{R}^N, \quad (50.11a)$$

$$\beta_h \|\mathbf{R}_\psi(\mathbf{P})\|_Q \leq \|\mathcal{B}^\top \mathbf{P}\|_{\ell_\varphi^2} \leq \|b\| \|\mathbf{R}_\psi(\mathbf{P})\|_Q, \quad \forall \mathbf{P} \in \mathbb{R}^M. \quad (50.11b)$$

*Proof.* See Exercise 50.4. □

### 50.2.2 Schur complement

Since the matrix  $\mathcal{A}$  is invertible, the vector  $\mathbf{U}$  can be eliminated from the linear system (50.9) yielding

$$\mathcal{S}\mathbf{P} = \mathcal{B}\mathcal{A}^{-1}\mathbf{F} - \mathbf{G}, \quad \mathcal{S} := \mathcal{B}\mathcal{A}^{-1}\mathcal{B}^\top. \quad (50.12)$$

Once  $\mathbf{P}$  is known,  $\mathbf{U}$  is obtained by solving  $\mathcal{A}\mathbf{U} = \mathbf{F} - \mathcal{B}^\top \mathbf{P}$ . The matrix  $\mathcal{S} \in \mathbb{R}^{M \times M}$  (up to a sign convention) is called *Schur complement* of  $\mathcal{A}$ ; see §49.2.1 for the infinite-dimensional counterpart. Notice that the matrix  $\mathcal{S}$  is invertible (if  $\mathcal{S}\mathbf{P} = \mathbf{0}$ , setting  $\mathbf{U} := -\mathcal{A}^{-1}\mathcal{B}^\top \mathbf{P}$ , we infer that  $\mathcal{C}(\mathbf{U}, \mathbf{P})^\top = (0, 0)^\top$ , and  $\mathcal{C}$  being invertible, this implies that  $\mathbf{U} = \mathbf{0}$  and  $\mathbf{P} = \mathbf{0}$ ).

Additional properties of the Schur complement matrix  $\mathcal{S}$  are available when the bilinear form  $a$  is symmetric and coercive, since in this case the matrix  $\mathcal{A}$  is symmetric positive definite.

**Proposition 50.13 (Symmetry and positivity of  $\mathcal{S}$ ).** *If  $\mathcal{A}$  is symmetric positive definite, so is  $\mathcal{S}$ .*

*Proof.* The definition of  $\mathcal{S}$  implies that  $\mathcal{S}^\top = \mathcal{B}(\mathcal{A}^{-1})^\top \mathcal{B}^\top$ , but  $(\mathcal{A}^{-1})^\top = (\mathcal{A}^\top)^{-1}$ . Hence,  $\mathcal{S}$  is symmetric if  $\mathcal{A}$  is symmetric. Let now  $\mathbf{P} \in \mathbb{R}^M$ . Then  $\mathbf{P}^\top \mathcal{S} \mathbf{P} = \mathbf{P}^\top \mathcal{B} \mathcal{A}^{-1} \mathcal{B}^\top \mathbf{P} = (\mathcal{B}^\top \mathbf{P})^\top \mathcal{A}^{-1} \mathcal{B}^\top \mathbf{P} \geq 0$ . This proves that  $\mathcal{S}$  is positive semidefinite. Moreover,  $\mathcal{S} \mathbf{P} = 0$  implies that  $\mathcal{B}^\top \mathbf{P} = 0$ , so that  $\mathbf{P} = 0$  since  $\mathcal{B}^\top$  has full column rank. Hence,  $\mathcal{S}$  is positive definite.  $\square$

Note that even if  $\mathcal{A}$  is symmetric positive definite, the matrix  $\mathcal{C}$  is symmetric but indefinite. Observing that

$$\mathcal{C} = \begin{pmatrix} \mathcal{I}_N & \mathbf{0}_{N,M} \\ \mathcal{B} \mathcal{A}^{-1} & \mathcal{I}_M \end{pmatrix} \begin{pmatrix} \mathcal{A} & \mathbf{0}_{N,M} \\ \mathbf{0}_{M,N} & -\mathcal{S} \end{pmatrix} \begin{pmatrix} \mathcal{I}_N & \mathcal{A}^{-1} \mathcal{B}^\top \\ \mathbf{0}_{N,M} & \mathcal{I}_M \end{pmatrix},$$

we infer from the Sylvester Law of Inertia (stating that two symmetric matrices  $\mathcal{C}$  and  $\mathcal{C}'$  satisfying  $\mathcal{C} = \mathcal{P} \mathcal{C}' \mathcal{P}^\top$  with  $\mathcal{P}$  invertible have the same number of positive, zero, and negative eigenvalues; see Golub and van Loan [218, p. 403]) that  $\mathcal{C}$  has  $N$  positive eigenvalues and  $M$  negative ones. Upper and lower bounds on the clusters of positive and negative eigenvalues of  $\mathcal{C}$  are derived in Rusten and Winther [338], Wathen and Silvester [390]. In practice, the matrix  $\mathcal{C}$  can be very poorly conditioned. We return to this issue in §50.3.2. Note that changing the lower-left block of  $\mathcal{C}$  into  $-\mathcal{B}$  produces a positive semidefinite, but nonsymmetric, matrix.

Let us now examine more closely the eigenvalues of  $\mathcal{S}$  (see Verfürth [375]). To this purpose, let  $\mathcal{M}_Q \in \mathbb{R}^{M \times M}$  be the matrix with entries  $\mathcal{M}_{Q,kl} := (\psi_k, \psi_l)_Q$  for all  $k, l \in \{1:M\}$ . The matrix  $\mathcal{M}_Q$  is symmetric by construction, and the identity  $\mathbf{P}^\top \mathcal{M}_Q \mathbf{P} = (\mathbf{R}_\psi(\mathbf{P}), \mathbf{R}_\psi(\mathbf{P}))_Q = \|\mathbf{R}_\psi(\mathbf{P})\|_Q^2$  for all  $\mathbf{P} \in \mathbb{R}^M$  shows that  $\mathcal{M}_Q$  is positive definite. Since  $Q$  is the  $L^2$ -space in many applications, the matrix  $\mathcal{M}_Q$  is called mass matrix (see §28.1.1). Let  $\mu_{\min}$  and  $\mu_{\max}$  be the lowest and largest eigenvalues of  $\mathcal{M}_Q$ . Recall from §28.2.1 that the (Euclidean) condition number  $\kappa(\mathcal{Z})$  of a symmetric invertible matrix  $\mathcal{Z}$  is the ratio of the largest to the smallest eigenvalues of  $\mathcal{Z}$  in absolute value.

**Proposition 50.14 (Spectrum of  $\mathcal{S}$ ).** *Assume that the bilinear form  $a$  is symmetric and coercive on  $V_h$  with constant  $\alpha_h$  and that the inf-sup condition (50.4b) for  $b$  is satisfied with constant  $\beta_h$ . Then the matrices  $\mathcal{S}$  and  $\mathcal{M}_Q$  are spectrally equivalent, i.e., the following holds true for all  $\mathbf{P} \in \mathbb{R}^M$ :*

$$\frac{\beta_h^2}{\|a\|} \leq \frac{\mathbf{P}^\top \mathcal{S} \mathbf{P}}{\mathbf{P}^\top \mathcal{M}_Q \mathbf{P}} \leq \frac{\|b\|^2}{\alpha_h}. \quad (50.13)$$

Moreover,  $\sigma(\mathcal{M}_Q^{-1} \mathcal{S}) \subset \left[ \frac{\beta_h^2}{\|a\|}, \frac{\|b\|^2}{\alpha_h} \right]$ , and  $\sigma(\mathcal{S}) \subset \left[ \mu_{\min} \frac{\beta_h^2}{\|a\|}, \mu_{\max} \frac{\|b\|^2}{\alpha_h} \right]$ , which implies that  $\kappa(\mathcal{M}_Q^{-1} \mathcal{S}) \leq \frac{\|a\|}{\alpha_h} \left( \frac{\|b\|}{\beta_h} \right)^2$  and  $\kappa(\mathcal{S}) \leq \frac{\|a\|}{\alpha_h} \left( \frac{\|b\|}{\beta_h} \right)^2 \kappa(\mathcal{M}_Q)$ .



*Proof.* (1) Proof of (50.13). For all  $P \in \mathbb{R}^M$ , we observe that

$$\begin{aligned} \sup_{Y \in \mathbb{R}^N} \frac{(\mathcal{B}^T P)^T Y}{(Y^T \mathcal{A} Y)^{\frac{1}{2}}} &= \sup_{Y \in \mathbb{R}^N} \frac{(\mathcal{B}^T P)^T \mathcal{A}^{-\frac{1}{2}} Y}{\|Y\|_{\ell^2(\mathbb{R}^N)}} = \sup_{Y \in \mathbb{R}^N} \frac{(\mathcal{A}^{-\frac{1}{2}} \mathcal{B}^T P)^T Y}{\|Y\|_{\ell^2(\mathbb{R}^N)}} \\ &= \|\mathcal{A}^{-\frac{1}{2}} \mathcal{B}^T P\|_{\ell^2(\mathbb{R}^N)} = (P^T \mathcal{S} P)^{\frac{1}{2}}, \end{aligned}$$

since  $\mathcal{A}$  is symmetric positive definite. Observing that  $\frac{1}{\|a\|} \leq \frac{\|R_\varphi(Y)\|_V^2}{Y^T \mathcal{A} Y} \leq \frac{1}{\alpha_h}$  for all  $Y \in \mathbb{R}^N$ , we infer that

$$\frac{1}{\|a\|} \|\mathcal{B}^T P\|_{\ell_\varphi^2}^2 \leq P^T \mathcal{S} P \leq \frac{1}{\alpha_h} \|\mathcal{B}^T P\|_{\ell_\varphi^2}^2.$$

Finally, (50.13) follows from (50.11b) using  $P^T \mathcal{M}_Q P = \|R_\psi(P)\|_Q^2$ .

(2) The spectrum and condition number for  $\mathcal{M}_Q^{-1} \mathcal{S}$  readily follow from (50.13), and the results for  $\mathcal{S}$  follow from the fact that  $\mu_{\min} \|P\|_{\ell^2(\mathbb{R}^M)}^2 \leq P^T \mathcal{M}_Q P \leq \mu_{\max} \|P\|_{\ell^2(\mathbb{R}^M)}^2$  for all  $P \in \mathbb{R}^M$ .  $\square$

### 50.2.3 Augmented Lagrangian for saddle point problems

Assume that the matrix  $\mathcal{A}$  is symmetric positive definite and that  $\mathcal{B}^T$  has full column rank. Referring to §49.3.2 for the infinite-dimensional setting we infer that the pair  $(U, P)$  solves the linear system (50.9) iff it is a saddle point of the Lagrangian

$$\mathcal{L}(Y, R) := \frac{1}{2} Y^T \mathcal{A} Y - F^T Y + R^T (\mathcal{B} Y - G). \quad (50.14)$$

Recall that

$$\inf_{Y \in \mathbb{R}^N} \sup_{R \in \mathbb{R}^M} \mathcal{L}(Y, R) = \mathcal{L}(U, P) = \sup_{R \in \mathbb{R}^M} \inf_{Y \in \mathbb{R}^N} \mathcal{L}(Y, R). \quad (50.15)$$

The optimization problem on the left-hand side of (50.15) amounts to minimizing the convex energy functional  $\mathfrak{E}(Y) := \frac{1}{2} Y^T \mathcal{A} Y - F^T Y$  over the affine subspace  $\{Y \in \mathbb{R}^N \mid \mathcal{B} Y = G\}$  since  $\sup_{R \in \mathbb{R}^M} \mathcal{L}(Y, R) = \infty$  if  $\mathcal{B} Y \neq G$ . Consider now the optimization problem on the right-hand side of (50.15). The minimization of  $\mathcal{L}(Y, R)$  over  $Y \in \mathbb{R}^N$  leads to the optimal solution  $Y_* := \mathcal{A}^{-1}(F - \mathcal{B}^T R)$ , and we are left with maximizing the following concave functional over  $\mathbb{R}^M$ :

$$R \mapsto \mathcal{L}(Y_*, R) = -\frac{1}{2} R^T \mathcal{S} R + (\mathcal{B} \mathcal{A}^{-1} F - G)^T R - \frac{1}{2} F^T \mathcal{A}^{-1} F,$$

where  $\mathcal{S} := \mathcal{B} \mathcal{A}^{-1} \mathcal{B}^T$ . The optimal solution to this maximization problem solves  $\mathcal{S} R = \mathcal{B} \mathcal{A}^{-1} F - G$ , i.e., we recover the Schur complement system (50.12).

The main idea of augmented Lagrangian methods (see Fortin and Glowinski [203]) is to add to the Lagrangian a least-squares penalty on the constraint. Specifically, letting  $\rho > 0$  be a real parameter and recalling the mass matrix  $\mathcal{M}_Q \in \mathbb{R}^{M \times M}$ , the *augmented Lagrangian* is defined as

$$\mathcal{L}_\rho(\mathbf{Y}, \mathbf{R}) := \mathcal{L}(\mathbf{Y}, \mathbf{R}) + \frac{\rho}{2}(\mathcal{B}\mathbf{Y} - \mathbf{G})^\top \mathcal{M}_Q^{-1}(\mathcal{B}\mathbf{Y} - \mathbf{G}).$$

Since we also have  $\mathcal{B}\mathbf{U} = \mathbf{G}$ , the solution to (50.9) is also the unique saddle point of the augmented Lagrangian  $\mathcal{L}_\rho$ , i.e.,  $(\mathbf{U}, \mathbf{P})$  can be found by solving the following linear system:

$$\begin{pmatrix} \mathcal{A}_\rho & \mathcal{B}^\top \\ \mathcal{B} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{U} \\ \mathbf{P} \end{pmatrix} = \begin{pmatrix} \mathbf{F}_\rho \\ \mathbf{G} \end{pmatrix}, \quad \begin{aligned} \mathcal{A}_\rho &:= \mathcal{A} + \rho \mathcal{B}^\top \mathcal{M}_Q^{-1} \mathcal{B}, \\ \mathbf{F}_\rho &:= \mathbf{F} + \rho \mathcal{B}^\top \mathcal{M}_Q^{-1} \mathbf{G}. \end{aligned} \quad (50.16)$$

The augmented Schur complement is defined as  $\mathcal{S}_\rho := \mathcal{B}\mathcal{A}_\rho^{-1}\mathcal{B}^\top$ . Recall that  $\sigma(\mathcal{M}_Q^{-1}\mathcal{S}) \subset [s_b, s_\sharp]$ , with  $s_b := \frac{\beta_h^2}{\|a\|}$  and  $s_\sharp := \frac{\|b\|^2}{\alpha_h}$ .

**Proposition 50.15 (Spectrum of  $\mathcal{S}_\rho$ ).** *The following holds true:*

$$\mathcal{S}_\rho^{-1} = \rho \mathcal{M}_Q^{-1} + \mathcal{S}^{-1}, \quad (50.17)$$

and  $\sigma(\mathcal{M}_Q^{-1}\mathcal{S}_\rho) \subset [(\rho + s_b^{-1})^{-1}, (\rho + s_\sharp^{-1})^{-1}]$  and  $\kappa(\mathcal{M}_Q^{-1}\mathcal{S}_\rho) \leq \frac{\rho + s_b^{-1}}{\rho + s_\sharp^{-1}}$ .

*Proof.* See Exercise 50.5 for the proof of (50.17). The properties on the spectrum and the condition number of  $\mathcal{S}_\rho$  follow readily.  $\square$

**Remark 50.16 (Value of  $\rho$ ).** Proposition 50.15 shows that taking  $\rho \gg 1$  improves the condition number of the Schur complement  $\mathcal{S}_\rho$ . A large value of  $\rho$  however deteriorates the conditioning of the matrix  $\mathcal{A}_\rho$  which makes it more difficult to invert iteratively. In practice, it is necessary to strike a balance between these two criteria.  $\square$

**Remark 50.17 (Hilbert setting).** The notion of augmented Lagrangian can be extended to the infinite-dimensional setting. The mass matrix  $\mathcal{M}_Q$  is then replaced by the Riesz–Fréchet isomorphism  $J_Q : Q \rightarrow Q'$ .  $\square$

The augmented Lagrangian technique is in general preferable to the following unconstrained penalty method:

$$\begin{pmatrix} \mathcal{A} & \mathcal{B}^\top \\ \mathcal{B} & -\epsilon \mathcal{M}_Q \end{pmatrix} \begin{pmatrix} \mathbf{U}_\epsilon \\ \mathbf{P}_\epsilon \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{G} \end{pmatrix}, \quad (50.18)$$

where  $\epsilon > 0$  is a small parameter. This technique is often referred to as artificial compressibility in the fluid mechanics literature. Eliminating  $\mathbf{P}_\epsilon$  from the first equation yields

$$\mathcal{A}_\frac{1}{\epsilon} \mathbf{U}_\epsilon = \mathbf{F} + \epsilon^{-1} \mathcal{B}^\top \mathcal{M}_Q^{-1} \mathbf{G}. \quad (50.19)$$

The advantage of (50.19) with respect to (50.9) (or to (50.16)) is that the system matrix (i.e.,  $\mathcal{A}_{\frac{1}{\epsilon}}$ ) is symmetric positive definite. The solution  $(\mathbf{U}_{\epsilon}, \mathbf{P}_{\epsilon})$  however differs from  $(\mathbf{U}, \mathbf{P})$ . In particular,  $\mathbf{U}$  fails to satisfy the constraint  $\mathcal{B}\mathbf{U} = \mathbf{G}$ , although the difference  $(\mathbf{U} - \mathbf{U}_{\epsilon}, \mathbf{P} - \mathbf{P}_{\epsilon})$  tends to zero as  $\epsilon \rightarrow 0$ . Unfortunately, taking  $\epsilon \ll 1$  makes the linear system (50.19) ill-conditioned.

**Proposition 50.18 (Penalty).** *Let  $\epsilon > 0$ . Let  $(\mathbf{U}, \mathbf{P})$  solve (50.9) and  $(\mathbf{U}_{\epsilon}, \mathbf{P}_{\epsilon})$  solve (50.18). The following holds true:*

$$\frac{\alpha_h \beta_h}{\|\mathbf{a}\|} \|\mathbf{R}_{\varphi}(\mathbf{U} - \mathbf{U}_{\epsilon})\|_V + \frac{\alpha_h \beta_h^2}{\|\mathbf{a}\|^2} \|\mathbf{R}_{\psi}(\mathbf{P} - \mathbf{P}_{\epsilon})\|_Q \leq \epsilon \|\mathbf{R}_{\psi}(\mathbf{P})\|_Q. \quad (50.20)$$

*Proof.* See Exercise 50.6. □

## 50.3 Iterative solvers

In this section, we discuss iterative solvers for the linear system (50.9) (or its augmented Lagrangian version (50.16)). First, we discuss the Uzawa algorithm as an example of a technique based on stationary iterations. Then, we present more efficient techniques based on preconditioned Krylov subspaces. We assume that the matrix  $\mathcal{A}$  is invertible and that the matrix  $\mathcal{B}^T$  has full column rank.

### 50.3.1 Uzawa algorithm

The Uzawa algorithm is an iterative method where  $\mathbf{U}$  and  $\mathbf{P}$  are updated one after the other. Given  $\mathbf{P}_0 \in \mathbb{R}^M$  and a parameter  $\eta > 0$ , the algorithm consists of constructing the iterates  $(\mathbf{U}_m, \mathbf{P}_m)$  for  $m = 1, 2, \dots$  as follows:

$$\mathcal{A}\mathbf{U}_m = \mathbf{F} - \mathcal{B}^T \mathbf{P}_{m-1}, \quad (50.21a)$$

$$\mathcal{M}_Q \mathbf{P}_m = \mathcal{M}_Q \mathbf{P}_{m-1} + \eta(\mathcal{B}\mathbf{U}_m - \mathbf{G}). \quad (50.21b)$$

This makes sense since  $\mathcal{A}$  and  $\mathcal{M}_Q$  are invertible. Eliminating  $\mathbf{U}_m$  gives

$$\mathcal{M}_Q \mathbf{P}_m = \mathcal{M}_Q \mathbf{P}_{m-1} - \eta(\mathcal{S}\mathbf{P}_{m-1} - \mathcal{B}\mathcal{A}^{-1}\mathbf{F} + \mathbf{G}). \quad (50.22)$$

In other words, the Uzawa algorithm is equivalent to the Richardson iteration applied to the linear system (50.12) left-preconditioned by the mass matrix  $\mathcal{M}_Q$ . (Recall that for a generic linear system  $\mathcal{Z}\mathbf{X} = \mathbf{Y}$ , the Richardson iteration reads  $\mathbf{X}_m = \mathbf{X}_{m-1} + \eta(\mathcal{Z}\mathbf{X}_{m-1} - \mathbf{Y})$ .) If  $\mathcal{A}$  is symmetric positive definite, we can use the bounds on the spectrum of  $\mathcal{M}_Q^{-1}\mathcal{S}$  from Proposition 50.14, that is,  $s_b := \frac{\beta_h^2}{\|\mathbf{a}\|} \leq \mathcal{M}_Q^{-1}\mathcal{S} \leq \frac{\|\mathbf{b}\|^2}{\alpha_h} := s_{\sharp}$  in the sense of quadratic forms. We then infer that the Richardson iteration (50.22) converges geometrically provided we take  $0 < \eta < \frac{2}{s_{\sharp}}$ , and the error reduction factor is maximized

by taking the optimal value  $\eta_{\text{opt}} := \frac{2}{s_b + s_{\sharp}}$ ; see Saad [339, p. 106]. It is often easier to estimate  $s_{\sharp}$  than  $s_b$  since  $\beta_h$  is more difficult to estimate than  $\alpha_h$ .

**Remark 50.19 (Implementation).** The matrices  $\mathcal{A}$  and  $\mathcal{B}$  are sparse (see §29.1), but  $\mathcal{S}$  is a dense matrix owing to the presence of  $\mathcal{A}^{-1}$  in the definition of  $\mathcal{S}$ . Since precomputing  $\mathcal{A}^{-1}$  is generally too expensive, an inner iteration has to be employed to compute the action of  $\mathcal{A}^{-1}$  on vectors in  $\mathbb{R}^N$ . The matrix  $\mathcal{B}$  can be assembled and stored once and for all, or its action on a given vector in  $\mathbb{R}^M$  can be computed on the fly whenever needed. In practice, one must often find a compromise between many (often conflicting) criteria: the memory space available; the number of times the linear system has to be solved; the ratio between the speed to access memory and the speed to perform floating point operations; parallelization; etc.  $\square$

**Remark 50.20 (Variants).** The mass matrix  $\mathcal{M}_Q$  can be replaced by the identity matrix  $\mathcal{I}_M$  in (50.21b). The advantage is that this avoids computing the inverse of the mass matrix (although this matrix is generally easy to invert since it is well-conditioned). The drawback is that the choice of the relaxation parameter  $\eta$  now depends on the spectrum of the unpreconditioned Schur complement matrix, which requires some information on the (mesh-dependent) spectrum of  $\mathcal{M}_Q$ . Another variant is to consider an approximate inverse of  $\mathcal{A}$  that is easy to compute, say  $\mathcal{H}$ , and to replace (50.21a) by  $\mathbf{U}_m = \mathbf{U}_{m-1} + \mathcal{H}(\mathbf{F} - \mathcal{A}\mathbf{U}_{m-1} - \mathcal{B}^T\mathbf{P}_{m-1})$  leading to an inexact Uzawa algorithm; see Bacuta [41] for a convergence analysis.  $\square$

### 50.3.2 Krylov subspace methods

Krylov subspace methods for solving (preconditioned) linear systems of the form (50.9) or variations thereof constitute an active area of research. In this section, we sketch a few important ideas and refer to Benzi et al. [52, §9] for a broader treatment and to Elman et al. [185, Chap. 6&8], Turek [366] for applications to fluid mechanics.

In the context of saddle point problems, the matrix  $\mathcal{C}$  in (50.9) is symmetric, but indefinite (recall that the matrix  $\mathcal{A}$  is symmetric positive definite by assumption). In this case, MINRES is a method of choice to solve (50.9); see [185, p. 289]. The attractive feature of MINRES is that it achieves an optimality property on the residual while employing only short-term recurrences. Specifically, at step  $m \geq 1$ , the iterate  $\mathbf{X}_m \in \mathbb{R}^{N+M}$  with residual  $\mathbf{R}_m := (\mathbf{F}, \mathbf{G})^T - \mathcal{C}\mathbf{X}_m$  satisfies the following optimality property (compare with Proposition 28.20 for the conjugate gradient method applied to symmetric positive definite linear systems):

$$\|\mathbf{R}_m\|_{\ell^2(\mathbb{R}^{N+M})} = \min_{\mathbf{Y} \in \mathbf{U}_0 + K_m} \|(\mathbf{F}, \mathbf{G})^T - \mathcal{C}\mathbf{Y}\|_{\ell^2(\mathbb{R}^{N+M})}, \quad (50.23)$$

with the Krylov subspace  $K_m := \text{span}\{\mathbf{R}_0, \mathcal{C}\mathbf{R}_0, \dots, \mathcal{C}^{m-1}\mathbf{R}_0\}$ . The convergence rate of MINRES depends on the spectrum of  $\mathcal{C}$ . More precisely, defining

$\tilde{c}_m := \min_{p \in \mathbb{P}_m, p(0)=1} \max_{\lambda \in \sigma(\mathcal{C})} |p(\lambda)|$ , one can prove that  $\|\mathbf{R}_m\|_{\ell^2(\mathbb{R}^{N+M})} \leq \tilde{c}_m \|\mathbf{R}_0\|_{\ell^2(\mathbb{R}^{N+M})}$  (this bound is sharp). The constant  $\tilde{c}_m$  can be estimated under the assumption that  $\sigma(\mathcal{C}) \subset [-a, -b] \cup [c, d]$  with positive real numbers  $a, b, c, d$  such that the two intervals have the same length (i.e.,  $d - c = a - b$ ). One can show that (see Greenbaum [221, Chap. 3])

$$\|\mathbf{R}_{2m}\|_{\ell^2(\mathbb{R}^{N+M})} \leq 2 \left( \frac{\sqrt{ad} - \sqrt{bc}}{\sqrt{ad} + \sqrt{bc}} \right)^m \|\mathbf{R}_0\|_{\ell^2(\mathbb{R}^{N+M})}. \quad (50.24)$$

The minimization property of MINRES implies that  $\|\mathbf{R}_{2m+1}\|_{\ell^2(\mathbb{R}^{N+M})} \leq \|\mathbf{R}_{2m}\|_{\ell^2(\mathbb{R}^{N+M})}$ , but it is possible that no reduction of the norm of the residual occurs in every other step, leading to a staircasing behavior of the iterates. A comparison of (28.23) with (50.24) shows that MINRES requires twice as many iterations as the Conjugate Gradient to reach a given threshold for a symmetric positive definite matrix with condition number  $\kappa^2$ . Hence, solving linear systems like (50.9) is a significant computational challenge, and preconditioning is essential. Before addressing this question let us observe that MINRES is bound to fail if  $\mathcal{A}$  is not symmetric, since symmetry is essential for MINRES to work properly. This happens, for instance, in fluid mechanics when solving the Oseen or (linearized) Navier–Stokes equations. One alternative is to use the GMRES method which retains an optimality property over the Krylov subspace at the price of storing a complete basis thereof; see Saad [339, §6.5] for a thorough description.

*Preconditioning* is a very important ingredient of Krylov subspace methods, especially for linear systems of the form (50.9). Here, we only discuss block preconditioners and refer the reader to [52, §10] and references therein for further insight into preconditioned Krylov methods. Block diagonal and triangular preconditioners are, respectively, of the form

$$\mathcal{P}_d := \begin{pmatrix} \hat{\mathcal{A}} & \mathbf{0} \\ \mathbf{0} & \hat{\mathcal{S}} \end{pmatrix}, \quad \mathcal{P}_t := \begin{pmatrix} \hat{\mathcal{A}} & \mathcal{B}^\top \\ \mathbf{0} & \hat{\mathcal{S}} \end{pmatrix}, \quad (50.25)$$

where  $\hat{\mathcal{A}}$  and  $\hat{\mathcal{S}}$  are easy-to-invert approximations of  $\mathcal{A}$  and  $\mathcal{S}$ . In the ideal case where  $\hat{\mathcal{A}} := \mathcal{A}$  and  $\hat{\mathcal{S}} := \mathcal{S}$ , a direct calculation shows that the left-preconditioned matrices  $\mathcal{P}_d^{-1}\mathcal{C}$  and  $\mathcal{P}_t^{-1}\mathcal{C}$  are zeroes of the polynomials  $p_d(\lambda) := (\lambda - 1)(\lambda - \frac{1 \pm \sqrt{5}}{2})$  and  $p_t(\lambda) := \lambda^2 - 1$ , respectively (see Kuznetsov [272], Murphy et al. [309]; see also (49.22) and Theorem 49.6), implying convergence in at most three (resp., two) steps for every preconditioned Krylov subspace method. The block triangular preconditioner  $\mathcal{P}_t$  breaks the symmetry of the system even if  $\mathcal{A}$  is symmetric, but this preconditioner is still quite effective in many cases, particularly for Oseen and Navier–Stokes flows (where the convective term breaks the symmetry of  $\mathcal{A}$  anyway). Note also that the costs of the two preconditioners in (50.25) are essentially identical since the cost of the additional multiplication by  $\mathcal{B}^\top$  is often marginal.

Effective choices for  $\hat{\mathcal{A}}$  and  $\hat{\mathcal{S}}$  are often driven by the application at hand. For Darcy's and Maxwell's equations (see Examples in §49.1.1 and §49.1.3),  $\mathcal{A}$  represents a zeroth-order differential operator (multiplication by a material property), and choosing a diagonal lumping for  $\hat{\mathcal{A}}$  together with some multilevel technique for  $\hat{\mathcal{S}}$  often works well if the material coefficients are smooth (see, e.g., Perugia and Simoncini [324] for magnetostatics problems). For the Stokes equations (see §49.1.2),  $\mathcal{A}$  represents a second-order differential operator, and the preconditioner  $\hat{\mathcal{A}}$  is typically based on some multilevel technique. The mass matrix associated with  $p$  can be used for  $\hat{\mathcal{S}}$  and a detailed eigenvalue analysis of the resulting block-diagonal preconditioned system can be found in Silvester and Wathen [347]. The approximation of the Schur complement becomes more delicate in the unsteady case and in the presence of convection. Preconditioners devised from the structure of the steady Navier–Stokes equations can be found in Elman et al. [185, Chap. 8] and the references therein. Furthermore, an attractive idea for transient and high-Reynolds number flows is to consider a block triangular preconditioner based on the augmented Lagrangian formulation (50.16) for the  $(1, 1)$ -block, together with the (scaled) mass matrix for the  $(2, 2)$ -block (thereby avoiding to consider the Schur complement); see Benzi and Olshanskii [51], Benzi et al. [53].

## Exercises

**Exercise 50.1 (Algebraic setting).** Let  $\mathcal{A} := \begin{pmatrix} 1 & \sqrt{2} \\ \sqrt{2} & 0 \end{pmatrix}$  and  $\mathcal{B} := (1, 0)^\top$ . Show that

$$\inf_{\mathbf{W} \in \ker(\mathcal{B})} \sup_{\mathbf{V} \in \ker(\mathcal{B})} \frac{\mathbf{W}^\top \mathcal{A} \mathbf{V}}{\|\mathbf{W}\|_{\ell^2(\mathbb{R}^2)} \|\mathbf{V}\|_{\ell^2(\mathbb{R}^2)}} < \inf_{\mathbf{V} \in \mathbb{R}^2} \sup_{\mathbf{W} \in \mathbb{R}^2} \frac{\mathbf{W}^\top \mathcal{A} \mathbf{V}}{\|\mathbf{W}\|_{\ell^2(\mathbb{R}^2)} \|\mathbf{V}\|_{\ell^2(\mathbb{R}^2)}}.$$

(*Hint*: one number is equal to 0 and the other is equal to 1.)

**Exercise 50.2 (Saddle point problem).** Let  $V, Q$  be Hilbert spaces and let  $a$  be a symmetric, coercive, bilinear form. Consider the discrete problem (50.2) and the bilinear form  $t(y, z) := a(v, w) + b(w, q) + b(v, r)$  for all  $y := (v, q), z := (w, r) \in X := V \times Q$ . Let  $X_h := V_h \times Q_h$  and consider the linear map  $P_h \in \mathcal{L}(X; X_h)$  such that for all  $x \in X$ ,  $P_h(x) \in X_h$  is the unique solution of  $t(P_h(x), y_h) = t(x, y_h)$  for all  $y_h \in X_h$ . Equip  $X$  and  $X_h$  with the norm  $\|(v, q)\|_{\tilde{X}} := (\|v\|_a^2 + \|q\|_Q^2)^{\frac{1}{2}}$  with  $\|v\|_a^2 := a(v, v)$ . (i) Prove that  $\|P_h\|_{\mathcal{L}(X; X)} \leq \tilde{c}_h := \frac{(4 \frac{\|b\|^2}{\alpha} + 1)^{\frac{1}{2}} + 1}{(4 \frac{\beta^2}{\|a\|} + 1)^{\frac{1}{2}} - 1}$ . (*Hint*: use Proposition 49.8.) (ii) Prove that  $\|u - u_h\|_a^2 + \|p - p_h\|_Q^2 \leq \tilde{c}_h^2 (\inf_{v_h \in V_h} \|u - u_h\|_a^2 + \inf_{q_h \in Q_h} \|p - q_h\|_Q^2)$ . (*Hint*: see the proof of Theorem 5.14.)

**Exercise 50.3 (Error estimate).** (i) Prove directly the estimate (50.7a) with  $c'_{1h}$  replaced by  $c''_{1h} := (1 + \frac{\|a\|}{\alpha_h})(1 + \frac{\|b\|}{\beta_h})$ . (*Hint:* consider  $z_h \in V_h$  s.t.  $B_h(z_h) := B_h(u_h - v_h)$  with  $v_h \in V_h$  arbitrary.) (ii) Assume that  $V$  is a Hilbert space,  $\ker(B_h) \subset \ker(B)$ , and  $g := 0$ . Prove that  $\|u - u_h\|_V \leq \frac{\|a\|}{\alpha_h} \inf_{v_h \in \ker(B_h)} \|u - v_h\|_V$ .

**Exercise 50.4 (Bound on  $\mathcal{A}$  and  $\mathcal{B}$ ).** (i) Prove Proposition 50.12. (*Hint:* observe that  $(\mathcal{A}U)^\top Y = a(\mathcal{R}_\varphi(U), \mathcal{R}_\varphi(Y))$ .) (ii) Let  $\mathcal{J}_V \in \mathbb{R}^{N \times N}$  be the symmetric positive definite matrix with entries  $\mathcal{J}_{V,ij} := (\varphi_i, \varphi_j)_X$  for all  $i, j \in \{1:N\}$ . Let  $\|\cdot\|_{\ell^2(\mathbb{R}^N)}$  denote the Euclidean norm in  $\mathbb{R}^N$ . Verify that  $\|\mathcal{R}_\varphi(U)\|_V = \|\mathcal{J}_V^{\frac{1}{2}}U\|_{\ell^2(\mathbb{R}^N)}$  and  $\|U\|_{\ell^2(\mathbb{R}^N)} = \|\mathcal{J}_V^{-\frac{1}{2}}U\|_{\ell^2(\mathbb{R}^N)}$  for all  $U \in \mathbb{R}^N$ .

**Exercise 50.5 ( $\mathcal{S}_\rho$ ).** The goal is to prove the identity (50.17). (i) Verify that  $\mathcal{A}_\rho^{-1} = \mathcal{A}^{-1} - \rho \mathcal{A}^{-1} \mathcal{B}^\top (\mathcal{M}_Q + \rho \mathcal{S})^{-1} \mathcal{B} \mathcal{A}^{-1}$ . (*Hint:* multiply the right-hand side by  $\mathcal{A}_\rho$  and develop the product.) (ii) Infer that  $\mathcal{S}_\rho = \mathcal{S} - \rho \mathcal{S} (\mathcal{M}_Q + \rho \mathcal{S})^{-1} \mathcal{S}$ . (iii) Conclude. (*Hint:* multiply the right-hand side by  $\rho \mathcal{M}_Q^{-1} + \mathcal{S}^{-1}$ .)

**Exercise 50.6 (Penalty).** (i) Prove Proposition 50.18. (*Hint:* verify that  $\mathcal{C}(U - U_\epsilon, P - P_\epsilon)^\top = (0, -\epsilon \mathcal{M}_Q P_\epsilon)^\top$  and use Proposition 50.12.) (ii) Replace the mass matrix  $\mathcal{M}_Q$  by the identity matrix  $\mathcal{I}_M$  times a positive coefficient  $\lambda$  in (50.18). Does the method still converge? Is there any interest of doing so? Can you think of another choice?

**Exercise 50.7 (Inexact Minres and DPG).** Let  $V, Y$  be Hilbert spaces and  $B \in \mathcal{L}(V; Y')$  be s.t.  $\beta \|v\|_V \leq \|B(v)\|_{Y'} \leq \|b\| \|v\|_V$  for all  $v \in V$  with  $0 < \beta \leq \|b\| < \infty$ . Set  $b(v, y) := \langle B(v), y \rangle_{Y', Y}$ . Let  $f \in Y'$ . Let  $J_Y : Y \rightarrow Y'$  denote the isometric Riesz–Fréchet isomorphism. (i) Show that the MINRES problem  $\min_{v \in V} \|f - B(v)\|_{Y'}$  has a unique solution  $u \in V$ . (*Hint:* introduce the sesquilinear form  $a(v, w) := \langle B(v), J_Y^{-1}(B(w)) \rangle_{Y', Y}$  and invoke the Lax–Milgram Lemma.) (ii) Let  $\{V_h \subset V\}_{h \in \mathcal{H}}$  and  $\{Y_h \subset Y\}_{h \in \mathcal{H}}$  be sequences of subspaces approximating  $V$  and  $Y$ , respectively. Assume that there is  $\beta_0 > 0$  s.t. for all  $h \in \mathcal{H}$ ,

$$\inf_{v_h \in V_h} \sup_{y_h \in Y_h} \frac{|b(v_h, y_h)|}{\|v_h\|_V \|y_h\|_{Y'}} \geq \beta_0. \quad (50.26)$$

Let  $I_h : Y_h \rightarrow Y$  be the canonical injection and  $I_h^* : Y' \rightarrow Y'_h$ . Show that the inexact MINRES problem  $\min_{v_h \in V_h} \|I_h^*(f - B(v_h))\|_{Y'_h}$  has a unique solution  $u_h \in V_h$ . (*Hint:* introduce the residual representative  $r_h := J_{Y_h}^{-1} I_h^*(f - B(u_h)) \in V_h$  and show that the pair  $(u_h, r_h) \in V_h \times Y_h$  solves a saddle point problem.) (iii) Show that the residual representative  $r_h \in Y_h$  is the unique solution of the following constrained minimization problem:  $\min_{z_h \in Y_h \cap (I_h^*(B(V_h)))^\perp} \frac{1}{2} \|z_h\|_{Y'}^2 - \langle I_h^*(f), z_h \rangle_{Y'_h, Y_h}$ . (*Hint:* see Proposition 49.11.) (iv) Assume now that  $f \in \text{im}(B)$  so that  $B(u) = f$ . Prove that there is  $c$  s.t.  $\|u - u_h\|_V \leq c \inf_{w_h \in V_h} \|u - w_h\|_V$  for all  $h \in \mathcal{H}$ . (*Hint:* use a Fortin operator.) *Note:* since  $\beta \|v_h\|_V \leq \|B(v_h)\|_{Y'}$  for all  $v_h \in V_h$ , it is natural to expect that the inf-sup condition (50.26) is satisfied if the subspace  $Y_h \subset Y$  is chosen rich

enough. The inexact residual minimization in a discrete dual norm is at the heart of the discontinuous Petrov–Galerkin (dPG) method; see Demkowicz and Gopalakrishnan [158], Gopalakrishnan and Qiu [219], Carstensen et al. [111]. The extension to reflexive Banach spaces is studied in Muga and van der Zee [308].