

Chapter 5

Metric Projection: applications

In this chapter we shall list a problems and applications which may be analyzed by the results of the last chapters and attacked by the methods presented in the next chapter. The applications are located in approximation theory, optimization, math finance, signal and image analysis and mathematical physics. Many of the problems may be considered as inverse problems; see [6, 33, 45].

5.1 Approximation problems

The first approximation problem considered in this section is the classical Chebyshev approximation problem. Then we consider the approximation by partial sums of a representation system in Hilbert and Banach spaces. Finally we mention a result which is the basis for discretization methods of elliptic equations.

Chebyshev approximation

A polynomial p is a function which can be written in the form

$$p(t) = a_0 + a_1 t + \cdots + a_m t^m$$

for some coefficients $a_0, \dots, a_m \in \mathbb{R}$. If $a_m \neq 0$, then the polynomial is said to be of degree m . Polynomials are just about the simplest mathematical functions that exist, requiring only multiplications and additions for their evaluation. Yet they also have the ability to represent very general nonlinear relationships. Approximation of more complicated functions by polynomials is a basic building block for a great many numerical techniques.

Polynomials are used to approximate a difficult to evaluate function, such as an interest function, a density or a distribution function, with the aim of fast evaluation on a computer. Here there is no interest in the polynomial curve itself. Rather the interest is on how closely the polynomial can follow the special function, and especially on how small the maximum error can be made. Very high order polynomials may be used here if they provide accurate approximations. Very often a function is not approximated directly, but is first transformed or standardized so to make it more amenable to polynomial approximation.

Here we sketch the best approximation of functions by polynomials in space $C[\mathbf{a}, \mathbf{b}]$ of continuous functions on the interval $[\mathbf{a}, \mathbf{b}]$. The distance we use is the supremum-norm

in $C[a, b]$:

$$\|f - g\|_\infty := \sup_{t \in [a, b]} |f(t) - g(t)|, \quad f, g \in C[a, b].$$

Let \mathcal{P}_n be the space of polynomials of degree less than or equal to n . The best approximation problem in the supremum norm (Chebyshev approximation problem) is the following one:

$$\text{Given } f \in \mathcal{P}_n \text{ find } p_0 \in \mathcal{P}_n \text{ such that } \|f - p_0\|_\infty = \text{dist}(f, \mathcal{P}_n). \quad (5.1)$$

This problem has a solution since \mathcal{P}_n is proximal in $C[a, b]$. The proximality is a consequence of the fact that for each $f \in C[a, b]$ the finite-dimensional space $\text{span}(\{f\} \cup \mathcal{P}_n)$ is reflexive.

Theorem 5.1 (Kolmogorov, 1948). *Let $f \in C[a, b]$ and $p_0 \in \mathcal{P}_n$. Then the following statements are equivalent:*

- (a) p_0 is a best approximation of f .
- (b) $\max_{t \in A_0} (f(t) - p_0(t))g(t) \geq 0$ for all $g \in \mathcal{P}_n$ where $A_0 := \{t \in [a, b] : |f(t) - p_0(t)| = \|f - p_0\|_\infty\}$ is the set of the so called **peak points**.

Notice that A_0 is nonempty and compact.

Proof:

Obviously, A_0 is compact and nonempty due to the continuity of $f - p_0$ and the compactness of $[a, b]$.

(a) \implies (b) Let $g \in \mathcal{P}_n$ and $\alpha > 0$. Then there exists $t_\alpha \in [a, b]$ with $\|f - p_0 + \alpha g\|_\infty \geq \|f - p_0\|_\infty$. Then

$$\begin{aligned} |f(t_\alpha) - p_0(t_\alpha) + \alpha g(t_\alpha)| &= \|f - p_0 + \alpha g\|_\infty \geq \|f - p_0\|_\infty \geq |f(t_\alpha) - p_0(t_\alpha)| \\ |f(t_\alpha) - p_0(t_\alpha)|^2 + 2(f(t_\alpha) - p_0(t_\alpha))g(t_\alpha) + \alpha^2|g(t_\alpha)|^2 &\geq |f(t_\alpha) - p_0(t_\alpha)|^2 \end{aligned}$$

This implies

$$(f(t_\alpha) - p_0(t_\alpha))g(t_\alpha) \geq -\frac{1}{2}\alpha^2\|g\|_\infty^2 \quad (5.2)$$

Taking the limit $\alpha \rightarrow 0$ we obtain a cluster point $t_0 \in [a, b]$ of the family t_α , $\alpha > 0$, due to the compactness of $[a, b]$. Then we obtain from the inequalities above $|f(t_0) - p_0(t_0)| \geq \|f - p_0\|_\infty$, i.e. $t_0 \in A_0$ and $(f(t_0) - p_0(t_0))g(t_0) \geq 0$.

(b) \implies (a) Let $g_1 \in \mathcal{P}_n$. Then we obtain $t_0 \in A_0$ with

$$(f(t_0) - p_0(t_0))(p_0(t_0) - g_1(t_0)) \geq 0$$

Since t_0 is a peak point of $f - p_0$ we have $|f(t_0) - p_0(t_0)| = \|f - p_0\|_\infty$ and we obtain

$$\begin{aligned} \|f - g_1\|_\infty^2 &\geq |f(t_0) - g_1(t_0)|^2 = |(f(t_0) - p_0(t_0)) + (p_0(t_0) - g_1(t_0))|^2 \\ &= |f(t_0) - p_0(t_0)|^2 + 2(f(t_0) - p_0(t_0))(p_0(t_0) - g_1(t_0)) + |p_0(t_0) - g_1(t_0)|^2 \\ &\geq |f(t_0) - p_0(t_0)|^2 = \|f - p_0\|_\infty^2. \end{aligned}$$

Since $g_1 \in \mathcal{P}_n$ was arbitrary chosen p_0 is a best approximation of f . ■

The characterization of a best approximation in Theorem 5.1 is proved on a direct way. Instead, we could refer to Theorem ??, but then we had to analyze the duality mapping in the space $C[a, b]$; see Example 3.11. This road to the result 5.1 can be found in [31].

Now, we want to derive the uniqueness of the best approximation. Here, an important role plays the so called Haar condition. Let $\mathbf{U} \subset C[\mathbf{a}, \mathbf{b}]$ be a subspace with $\dim(\mathbf{U}) = n + 1$. We say that \mathbf{U} satisfies the **Haar condition** if and only the property

$$\text{Each } p \in \mathcal{P}_n, p \neq \theta, \text{ has at most } n \text{ zeros} \quad (5.3)$$

holds. If the Haar condition holds we call \mathbf{U} a **Haar space**. Clearly the subspace $\mathbf{U} := \mathcal{P}_n$ satisfies the Haar condition due to the fundamental theorem of algebra.

Let t_0, t_1, \dots, t_n be pairwise different points in $[\mathbf{a}, \mathbf{b}]$. If we choose a basis ϕ_0, \dots, ϕ_n in a Haar space, then the matrix $B := (\phi_i(t_j))_{i,j=0,\dots,n}$ has a non-zero determinant. A consequence of this fact is that for $n + 1$ distinct points t_0, \dots, t_n in $[\mathbf{a}, \mathbf{b}]$ and data $y_0, \dots, y_n \in \mathbb{R}$ the interpolation problem

$$g(t_i) = y_i, i = 0, \dots, n, \quad (5.4)$$

has a uniquely determined solution g in the Haar space.

Lemma 5.2. *Let $p_0 \in \mathcal{P}_n$ be the best approximation of $f \in C[\mathbf{a}, \mathbf{b}] \setminus \mathcal{P}_n$. Then there exist at least $n + 1$ distinct points ξ with $|f(\xi) - p_0(\xi)| = \|f - p_0\|_\infty$.*

Proof:

In a first step we prove that there exist at least two distinct points $t_0, t_1 \in [\mathbf{a}, \mathbf{b}]$ with

$$f(t_0) - p_0(t_0) = -(f(t_1) - p_0(t_1)) = \|f - p_0\|_\infty = E := \text{dist}(f, \mathcal{P}_n).$$

Let $t_1 \in [0, 1]$ with $f(t_1) - p_0(t_1) = E$. If the assertion is false, then we have $e := \min_{\xi \in [\mathbf{a}, \mathbf{b}]} (f(\xi) - p_0(\xi)) > -E$. Hence $e + E \neq 0$ and the polynomial $q := p_0 + \frac{1}{2}(E + e)$ belongs to \mathcal{P}_n . We have

$$E - \frac{1}{2}(E + e) \geq f(t) - p_0(t) - \frac{1}{2}(E + e) \geq e - \frac{1}{2}(E + e) \text{ for all } t \in [0, 1]$$

which implies

$$-\frac{1}{2}(E - e) \leq f(t) - q(t) \leq \frac{1}{2}(E - e).$$

This shows

$$\|f - q\|_\infty \leq \frac{1}{2}(E - e) < E = \|f - p_0\|_\infty = \text{dist}(f, \mathcal{P}_n).$$

This is a contradiction.

Suppose that there exist only $m \leq n + 1$ distinct points $\xi_0, \dots, \xi_m \in [\mathbf{a}, \mathbf{b}]$ with $|f(\xi_i) - p_0(\xi_i)| = \|f - p_0\|_\infty, i = 0, \dots, m$. Then we can find a polynomial $q \in \mathcal{P}_n$ with $q(\xi_i) = -(f(\xi_i) - p_0(\xi_i)), i = 0, \dots, m$. This implies

$$\max_{i=0,\dots,m} (f(\xi_i) - p_0(\xi_i))q(\xi_i) = \max_{i=0,\dots,m} (-(f(\xi_i) - p_0(\xi_i))^2) < 0.$$

This is a contradiction to condition (b) in Theorem 5.1. ■

Theorem 5.3. \mathcal{P}_n is a Chebyshev subspace of $C[\mathbf{a}, \mathbf{b}]$.

Proof:

We know already that \mathcal{P}_n is proximal. Let $f \in C[\mathbf{a}, \mathbf{b}]$. If $f \in \mathcal{P}_n$ then the best approximation in \mathcal{P}_n is given by f and therefore uniquely determined. Let $f \in C[\mathbf{a}, \mathbf{b}] \setminus \mathcal{P}_n$. Suppose that $\mathbf{p}_0, \mathbf{p}_1 \in \mathcal{P}_n$ are best approximations of f . Then with $\mathbf{p} := \frac{1}{2}(\mathbf{p}_0 + \mathbf{p}_1)$ we have

$$\|f - \mathbf{p}_0\|_\infty = \|f - \mathbf{p}_1\|_\infty = \|f - \mathbf{p}\|_\infty = \text{dist}(f, \mathcal{P}_n) > 0.$$

Hence,

$$\text{dist}(f, \mathcal{P}_n) \leq \|f - \mathbf{p}\|_\infty \leq \frac{1}{2}\|f - \mathbf{p}_0\|_\infty + \frac{1}{2}\|f - \mathbf{p}_1\|_\infty = \text{dist}(f, \mathcal{P}_n).$$

This shows that \mathbf{p} is also a best approximation. By Lemma 5.2 there exist $n + 2$ distinct points $\mathbf{t}_0, \dots, \mathbf{t}_{n+1} \in [\mathbf{a}, \mathbf{b}]$ with $|f(\mathbf{t}_i) - \mathbf{p}(\mathbf{t}_i)| = \text{dist}(f, \mathcal{P}_n)$, $i = 0, \dots, n + 1$. For a fixed \mathbf{t}_k we have

$$f(\mathbf{t}_k) - \frac{1}{2}(\mathbf{p}_0(\mathbf{t}_k) + \mathbf{p}_1(\mathbf{t}_k)) = f(\mathbf{t}_k) - \mathbf{p}(\mathbf{t}_k) = \text{dist}(f, \mathcal{P}_n) = \|f - \mathbf{p}_0\|_\infty \geq f(\mathbf{t}_k) - \mathbf{p}_0(\mathbf{t}_k),$$

which gives $\mathbf{p}_1(\mathbf{t}_k) \leq \mathbf{p}_0(\mathbf{t}_k)$. Similarly,

$$f(\mathbf{t}_k) - \frac{1}{2}(\mathbf{p}_0(\mathbf{t}_k) + \mathbf{p}_1(\mathbf{t}_k)) = f(\mathbf{t}_k) - \mathbf{p}(\mathbf{t}_k) = \text{dist}(f, \mathcal{P}_n) = \|f - \mathbf{p}_1\|_\infty \geq f(\mathbf{t}_k) - \mathbf{p}_1(\mathbf{t}_k),$$

leads to $\mathbf{p}_0(\mathbf{t}_k) \leq \mathbf{p}_1(\mathbf{t}_k)$. Thus, we have $\mathbf{p}_0(\mathbf{t}_k) = \mathbf{p}_1(\mathbf{t}_k)$, $k = 0, \dots, n + 1$. On the other hand, if $f(\mathbf{t}_k) - \mathbf{p}(\mathbf{t}_k) = -\text{dist}(f, \mathcal{P}_n)$, then

$$f(\mathbf{t}_k) - \frac{1}{2}(\mathbf{p}_0(\mathbf{t}_k) + \mathbf{p}_1(\mathbf{t}_k)) = f(\mathbf{t}_k) - \mathbf{p}(\mathbf{t}_k) = -\text{dist}(f, \mathcal{P}_n) = -\|f - \mathbf{p}_0\|_\infty \leq f(\mathbf{t}_k) - \mathbf{p}_0(\mathbf{t}_k),$$

which gives $\mathbf{p}_0(\mathbf{t}_k) \leq \mathbf{p}_1(\mathbf{t}_k)$. A similar procedure gives then $\mathbf{p}_0(\mathbf{t}_k) \geq \mathbf{p}_1(\mathbf{t}_k)$. Hence, we have in this case also $\mathbf{p}_0(\mathbf{t}_k) = \mathbf{p}_1(\mathbf{t}_k)$, $k = 0, \dots, n + 1$. We conclude that $\mathbf{p}_0 - \mathbf{p}_1 \in \mathcal{P}_n$ has $n + 2$ distinct zeroes and hence $\mathbf{p}_0 = \mathbf{p}_1$. \blacksquare

5.2 Greedy approximation

Greedy approximation in Hilbert space

A generic problem of mathematical and numerical analysis is to approximately represent a given function. It is a classical problem that goes back to the first results on Taylor's and Fourier's expansions of a function. The first step in solving the representation problem is to choose a representation system. Traditionally, a representation system has natural features such as minimality, orthogonality, simple structure and nice computational characteristics. The most typical representation systems are the trigonometric system $(\mathbf{e}^{ikt})_{k \in \mathbb{Z}}$, the algebraic system $(\mathbf{t}^k)_{k \in \mathbb{N}}$, the spline system, the finite element system and the wavelet system and their various variants. In general we may speak of a basis sequence $(\mathbf{b}^k)_{k \in \mathbb{N}}$ in a Hilbert space or Banach space. The second step in solving the representation problem is to choose a form of an approximant that is built on the basis of the chosen representation system.

Let us consider the Hilbert space case. Let \mathcal{H} be a separable Hilbert space and let the representation system $(\mathbf{e}^k)_{k \in \mathbb{N}}$ be an orthonormal basis. As we know, every element $\mathbf{x} \in \mathcal{H}$ has a representation

$$\mathbf{x} = \sum_{i \in \mathbb{N}} \langle \mathbf{x} | \mathbf{e}^i \rangle \mathbf{e}^i.$$

In a classical way that was used for centuries, an approximant for an element $\mathbf{x} \in \mathcal{H}$ is a partial sum

$$\mathbf{S}_m(\mathbf{x}) := \sum_{i=1}^m \langle \mathbf{x} | \mathbf{e}^i \rangle \mathbf{e}^i \quad (m \in \mathbb{N})$$

One calls such a construction of an approximant a **linear m-term approximation**.

In order to put this procedure into the terminology of the last chapter, we set $\mathbf{U}_m := \text{span}(\{\mathbf{e}^1, \dots, \mathbf{e}^m\})$ and consider the approximation problem

$$\text{Find } \mathbf{w} \in \mathbf{U}_m \text{ which satisfies } \|\mathbf{x} - \mathbf{w}\| = \inf_{\mathbf{u} \in \mathbf{U}_m} \|\mathbf{x} - \mathbf{u}\| = \text{dist}(\mathbf{x}, \mathbf{U}_m).$$

We know that this problem has a uniquely determined solution which we denote by $\mathbf{P}_{\mathbf{U}_m}(\mathbf{x})$ again; $\mathbf{P}_{\mathbf{U}_m}$ is the associated metric projection. From the fact that $\mathbf{x} - \mathbf{S}_m(\mathbf{x})$ is orthogonal to \mathbf{U}_m we conclude that $\mathbf{P}_{\mathbf{U}_m}(\mathbf{x}) = \mathbf{S}_m(\mathbf{x})$ holds true; see (4) in Theorem 3.3.

Experience in signal and image processing has led to the suggestion that it is more beneficial to use a general m -term approximation with respect to the basis chosen. This means that for $\mathbf{x} \in \mathcal{H}$ we look for an approximation of the form

$$\mathbf{S}_{\Lambda_m}(\mathbf{x}) := \sum_{i \in \Lambda_m} \langle \mathbf{x} | \mathbf{e}^i \rangle \mathbf{e}^i$$

where Λ_m is a subset of m indices. The choice of this set Λ_m has to be chosen appropriate with respect to \mathbf{x} . Notice, both $\mathbf{S}_m(\mathbf{x})$ and $\mathbf{S}_{\Lambda_m}(\mathbf{x})$, are of the same complexity. From the Parseval equality we obtain

$$\|\mathbf{x} - \mathbf{S}_{\Lambda_m}(\mathbf{x})\|^2 = \left\| \sum_{i \in \mathbb{N}, i \notin \Lambda_m} \langle \mathbf{x} | \mathbf{e}^i \rangle \mathbf{e}^i \right\|^2 = \sum_{i \in \mathbb{N}, i \notin \Lambda_m} |\langle \mathbf{x} | \mathbf{e}^i \rangle|^2 =: \eta_{\Lambda_m}(\mathbf{x}).$$

This type of approximation is a nonlinear approach since for a fixed m the approximants come from different subspaces spanned by $\mathbf{e}^i, i \in \Lambda_m(\mathbf{x})$.

The key for an optimal choice of the set $\Lambda_m(\mathbf{x})$ is the simple observation that the approximation error $\eta_{\Lambda_m}(\mathbf{x})$ will be minimized if the set Λ_m contains those indices whose modulus of the coefficients is biggest. This choice is called the **greedy choice**. This means that the **greedy subset** Λ_m is chosen in such a way that

$$\min_{i \in \Lambda_m} |\langle \mathbf{x} | \mathbf{e}^i \rangle| \geq \max_{i \notin \Lambda_m} |\langle \mathbf{x} | \mathbf{e}^i \rangle|$$

Clearly, a greedy subset is not necessarily uniquely determined. We set

$$\sigma_m(\mathbf{x}) := \|\mathbf{x} - \mathbf{S}_{\Lambda_m}(\mathbf{x})\| \text{ where } \Lambda_m \text{ is a greedy subset.}$$

Since the approximation of \mathbf{x} by $\mathbf{S}_m(\mathbf{x})$ can be considered as a m -term approximation we obtain

$$\sigma_m(\mathbf{x}) \leq \text{dist}(\mathbf{x}, \mathbf{U}_m).$$

In order to confirm that greedy approximation may decrease the approximation error let $m \in \mathbb{N}$ and choose $\mathbf{x} := \mathbf{e}^{m+1}$. We obtain $\mathbf{S}_m(\mathbf{x}) = \mathbf{0}$ with $\|\mathbf{x} - \mathbf{S}_m(\mathbf{x})\| = 1$, whereas $\mathbf{S}_{\Lambda_m}(\mathbf{x}) = \mathbf{x}$ and no error occur.

Algorithm 5.1 Greedy algorithm in Hilbert space

Given a Hilbert space with an orthonormal basis $(\mathbf{e}^i)_{i \in \mathbb{N}}$. Let $\mathbf{x} \in \mathcal{H}$. This algorithm computes a greedy subset Λ_m and a greedy approximation $\mathbf{S}_{\Lambda_m}(\mathbf{x})$.

- (1) $\mathbf{x}_0 := \mathbf{x}, j := 0, \Lambda_0 := \emptyset$.
 - (2) For $j = 1, \dots, m$ repeat
 - (i) Find n_j such that $|\langle \mathbf{x}_{j-1} | \mathbf{e}^{n_j} \rangle| \geq |\langle \mathbf{x}_{j-1} | \mathbf{e}^k \rangle|$ for all $k \in \mathbb{N}$.
 - (ii) Set $\Lambda_j := \Lambda_{j-1} \cup \{n_j\}, \mathbf{x}_j := \mathbf{x}_{j-1} - \langle \mathbf{x}_{j-1} | \mathbf{e}^{n_j} \rangle \mathbf{e}^{n_j}$.
 - (3) Set $\mathbf{S}_{\Lambda_m} := \mathbf{x}_m$.
-

Greedy approximation in Banach space

In order to translate the considerations above to the case of Banach spaces we need the terminology of a basis in a Banach space. The appropriate concept is that of a Schauder basis. We follow mainly [48, 49] and [12].

Definition 5.4. Let \mathcal{X} be a Banach space. A sequence $(\mathbf{b}^i)_{i \in \mathbb{N}}$ is called a **Schauder basis** if for each $\mathbf{x} \in \mathcal{X}$ there exist uniquely determined coefficients $c_i, i \in \mathbb{N}$, such that $\mathbf{x} = \sum_{i=1}^{\infty} c_i \mathbf{b}^i$. \square

Necessarily, a Banach space with a Schauder basis is separable (choose rational coefficients for the approximating series). Therefore, \mathbf{l}_{∞} possesses no Schauder basis. In the spaces $\mathbf{l}_p, 1 \leq p < \infty$, the standard unit vectors build a Schauder basis. The **Haar system** is a Schauder basis in $L_p[0, 1], 1 \leq p < \infty$. It was for a long time an open problem whether each separable Banach space possesses a Schauder basis. A negative answer has been given by P. Enflo [22] in 1973 when he constructed a Banach space which satisfies not the approximation property since this property is a necessary condition for the existence of a Schauder basis. The approximation property may be formulated as follows:

A Banach space \mathcal{X} has the approximation property if the following holds:
 For each compact subset K of \mathcal{X} and for each $\varepsilon > 0$ there exists a linear bounded operator $T : \mathcal{X} \rightarrow \mathcal{X}$ with $\dim \text{ran}(T) < \infty$ and $\|T(\mathbf{x}) - \mathbf{x}\| \leq \varepsilon$ for all $\mathbf{x} \in K$.

Let \mathcal{X} be a separable Banach space with a Schauder basis $(\mathbf{b}^i)_{i \in \mathbb{N}}$. Again, we denote the partial sums by $\mathbf{S}_m(\mathbf{x})$:

$$\mathbf{S}_m(\mathbf{x}) = \sum_{i=1}^m c_i \mathbf{b}^i \text{ if } \mathbf{x} = \sum_{i=1}^{\infty} c_i \mathbf{b}^i.$$

Clearly, the mappings

$$\lambda_n : \mathcal{X} \ni x = \sum_{i=1}^{\infty} c_i b^i \longmapsto c_n \in \mathbb{R}, n \in \mathbb{N},$$

are well defined linear functionals due to the uniqueness of the presentation of x as a series. As a consequence, the mapping $x \longmapsto S_m(x)$ is linear too. Moreover, one can show that for each $n \in \mathbb{N}$ the mapping $x \longmapsto S_m(x)$ is a bounded linear operator and the number $BC := \sup_{m \in \mathbb{N}} \|S_m\|$ is finite. Now, the coefficient functionals λ_n are linear and bounded too.

We set $U_m := \text{span}(\{b^1, \dots, b^m\})$ and consider the approximation problem

$$\text{Find } w \in U_m \text{ which satisfies } \|x - w\| = \inf_{u \in U_m} \|x - u\| = \text{dist}(x, U_m).$$

Since U_m is finite dimensional this problem has a solution but uniqueness is not guaranteed. Moreover, it is not clear that $S_m(x)$ is a solution of the problem. Since

$$\|S_m(x)\| \leq BC \|x\| \text{ for all } m \in \mathbb{N}, x \in \mathcal{X}.$$

we have for $u \in U_m$

$$\|x - S_m(x)\| \leq \|x - u + u - S_m(x)\| \leq \|x - u\| + \|S_m(x - u)\| \leq (1 + BC)\|x - u\|.$$

Hence

$$\|x - S_m(x)\| \leq (1 + BC)\text{dist}(x, U_m). \quad (5.5)$$

(5.5) shows that $S_m(x)$ is almost a best approximation of x in U_m . Moreover, $\lim_m S_m(x) = x$ due to the definition of a Schauder basis.

Similar to the Hilbert space case, we can consider the greedy approximation in the Banach space \mathcal{X} with Schauder basis $(b^i)_{i \in \mathbb{N}}$. The error of the best greedy approximation is

$$\sigma_m(x) := \inf_{\Lambda \subset \mathbb{N}, \#\Lambda=m} \inf\{\|x - \sum_{i \in \Lambda} c_i b^i\| : c_i \in \mathbb{R}, i \in \Lambda\}$$

Obviously, $\sigma_m(x) \leq \text{dist}(x, U_m)$ for all $m \in \mathbb{N}, x \in \mathcal{X}$. Moreover, we have $\sigma_m(b_{m+1}) < \text{dist}(x, U_m)$.

A permutation $\rho : \mathbb{N} \longrightarrow \mathbb{N}$ is called **decreasing** for x if

$$|\lambda_{\rho(i)}(x)| \geq |\lambda_{\rho(i+1)}(x)| \text{ for all } i \in \mathbb{N}$$

holds true. Now, we can define the family of greedy approximants: for $x \in \mathcal{X}$ and a decreasing permutation ρ we call

$$G_m(x, \rho) := \sum_{i=1}^m \lambda_{\rho(i)}(x) b^{\rho(i)}$$

the **m-th greedy-approximant**. The m-th greedy approximant is a partial sum of a rearranged series expansion of x . Therefore, convergence of $G_m(x, \rho) \rightarrow x$ for all $x \in \mathcal{X}$ can be expected only if every rearrangement of the basis $(b^i)_{i \in \mathbb{N}}$ is again a Schauder basis.

Definition 5.5. Let \mathcal{X} be a Banach space. A sequence $(\mathbf{b}^i)_{i \in \mathbb{N}}$ is called a **unconditional Schauder basis** if for each permutation $\pi : \mathbb{N} \rightarrow \mathbb{N}$ $(\mathbf{b}^{\pi(i)})_{i \in \mathbb{N}}$ is a Schauder basis. \square

For instance, in [2, 19, 18, 35] one can find results under which additional assumptions convergence of the greedy-approximants holds. The main source of different behaviour of greedy-approximants are wavelet bases.

5.3 Cea's Lemma

Cea's Lemma is fundamental for error estimates in finite element methods. It describes the bridge between the discretization error by finite elements in a (partial) differential equation and the best approximation error by the chosen finite elements.

Let \mathcal{V} be a Hilbert space with inner product $\langle \cdot | \cdot \rangle$ and let $\mathbf{a} : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ be a mapping with the properties

- (1) \mathbf{a} is bilinear.
- (2) \mathbf{a} is continuous, i.e. there exists a constant $\mathbf{c}_1 > 0$ such that

$$|\mathbf{a}(\mathbf{u}, \mathbf{v})| \leq \mathbf{c}_1 \|\mathbf{u}\| \|\mathbf{v}\| \text{ for all } \mathbf{u}, \mathbf{v} \in \mathcal{V}.$$

- (3) \mathbf{a} is coercive, i.e. there exist a constant \mathbf{c}_2 such that

$$\mathbf{a}(\mathbf{u}, \mathbf{u}) \geq \mathbf{c}_2 \|\mathbf{u}\|^2 \text{ for all } \mathbf{u} \in \mathcal{V}.$$

Consider for $\mathbf{y} \in \mathcal{V}$ the problem

$$\text{Find } \mathbf{u} \in \mathcal{V} \text{ such that } \mathbf{a}(\mathbf{u}, \mathbf{v}) = \langle \mathbf{y} | \mathbf{v} \rangle \text{ for all } \mathbf{v} \in \mathcal{V}. \quad (5.6)$$

To solve this problem one chooses a finite-dimensional subspace \mathbf{V}_h of \mathcal{V} (the space of „finite elements“) and considers the discretized problem

$$\text{Find } \mathbf{u}_h \in \mathbf{V}_h \text{ such that } \mathbf{a}(\mathbf{u}_h, \mathbf{v}) = \langle \mathbf{y} | \mathbf{v} \rangle \text{ for all } \mathbf{v} \in \mathbf{V}_h. \quad (5.7)$$

The parameter \mathbf{h} is measure which describes the meshwidth of the finite elements and determines the dimension of \mathbf{V}_h . By the Lax-Milgram lemma, each of the problems above has exactly one solution; see for instance [51].

Lemma 5.6 (Cea's Lemma). *Let \mathbf{u}, \mathbf{u}_h be the solution of (5.6) and (5.8), respectively. Under the assumptions above we have*

$$\|\mathbf{u} - \mathbf{u}_h\| \leq \mathbf{c}_1 (\mathbf{c}_2)^{-1} \text{dist}(\mathbf{u}, \mathbf{V}_h). \quad (5.8)$$

Proof:

We have for all $\mathbf{v} \in \mathbf{V}_h$

$$\mathbf{c}_2 \|\mathbf{u} - \mathbf{u}_h\|^2 \leq \mathbf{a}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) = \mathbf{a}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{v}) + \mathbf{a}(\mathbf{u} - \mathbf{u}_h, \mathbf{v} - \mathbf{u}_h).$$

Since

$$\mathbf{a}(\mathbf{u}, \mathbf{v}) = \langle \mathbf{y} | \mathbf{v} \rangle = \mathbf{a}(\mathbf{u}_h, \mathbf{v}) \text{ for all } \mathbf{v} \in \mathcal{V}_h.$$

we obtain

$$c_2 \|\mathbf{u} - \mathbf{u}_h\|^2 \leq \mathbf{a}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{v}) \leq c_1 \|\mathbf{u} - \mathbf{u}_h\| \|\mathbf{u} - \mathbf{v}\| \text{ for all } \mathbf{v} \in \mathcal{V}_h$$

and the result follows. ■

Lemma 5.6 says that \mathbf{u}_h is the „best approximation“ of \mathbf{u} in \mathcal{V}_h , up to a constant $c_1(c_2)^{-1}$. If the bilinear form \mathbf{a} is in addition symmetric, i.e. $\mathbf{a}(\mathbf{v}, \mathbf{w}) = \mathbf{a}(\mathbf{w}, \mathbf{v})$ for all $\mathbf{v}, \mathbf{w} \in \mathcal{V}$, then the constant $c_1(c_2)^{-1}$ may be replaced by the constant $\sqrt{c_1(c_2)^{-1}}$.

Example 5.7. Consider the boundary value problem

$$-\mathbf{u}'' + \mathbf{p}\mathbf{u}' + \mathbf{u} = \mathbf{f}, \quad \xi \in (0, 1), \quad \mathbf{u}(0) = \mathbf{u}(1) = 0, \quad (5.9)$$

where \mathbf{f} is a given function and \mathbf{p} is a constant. The problem consists in finding a function \mathbf{u} which satisfies the differential equation and the boundary conditions in (5.9).

If multiply the equation by a function

$$\mathbf{v} \in C_0^1(0, 1) := \{\mathbf{w} : (0, 1) \rightarrow \mathbb{R} : \mathbf{w} \text{ differentiable and } \mathbf{w}(0) = \mathbf{w}(1) = 0\}$$

then we obtain by formal computational steps

$$\mathbf{a}(\mathbf{u}, \mathbf{w}) = \int_0^1 \mathbf{u}'\mathbf{w}' + \mathbf{p}\mathbf{u}'\mathbf{w} + \mathbf{u}\mathbf{w} \, d\xi = \int_0^1 \mathbf{f}\mathbf{w} \, d\xi.$$

In order to obtain a Hilbert space setting we have to consider the completion of $C_0^2(0, 1)$ under the norm

$$\|\mathbf{w}\|_1 := \left(\int_0^1 (\mathbf{w}'^2 + \mathbf{w}^2) \, d\xi \right)^{\frac{1}{2}} = (\|\mathbf{w}'\|_0^2 + \|\mathbf{w}\|_0^2)^{\frac{1}{2}}$$

where $\|\cdot\|_0$ is the norm in $L_2(0, 1)$ getting the Hilbert space $\mathcal{V} := H_0^1(0, 1)$; $H_0^1(0, 1)$ is endowed with the inner product

$$\langle \mathbf{v}, \mathbf{w} \rangle := \int_0^1 (\mathbf{v}'\mathbf{w}' + \mathbf{v}\mathbf{w}) \, d\xi, \quad \mathbf{v}, \mathbf{w} \in H_0^1(0, 1).$$

\mathbf{a} is then defined on $\mathcal{V} \times \mathcal{V}$ as above. \mathbf{a} is symmetric if $\mathbf{p} = 0$, otherwise non-symmetric. For convenience, assume $0 \leq \mathbf{p} < 2$. The continuity of \mathbf{a} follows from

$$|\mathbf{a}(\mathbf{u}, \mathbf{v})| \leq |\langle \mathbf{u} | \mathbf{v} \rangle| + \mathbf{p} \left| \int_0^1 \mathbf{u}'\mathbf{v} \, d\xi \right| \leq \|\mathbf{u}\|_1 \|\mathbf{v}\|_1 + \mathbf{p} \|\mathbf{u}'\|_0 \|\mathbf{v}\|_0 \leq 3 \|\mathbf{u}\|_1 \|\mathbf{v}\|_1$$

The coercivity of \mathbf{a} is implied by

$$\begin{aligned} \mathbf{a}(\mathbf{v}, \mathbf{v}) &= \int_0^1 (\mathbf{v}'^2 + \mathbf{p}\mathbf{v}'\mathbf{v} + \mathbf{v}^2) \, d\xi \\ &\geq \frac{1}{2} \int_0^1 (\mathbf{v}'^2(2 - \mathbf{p}) + \mathbf{v}^2(2 - \mathbf{p})) \, d\xi \\ &\geq \frac{1}{2} (2 - \mathbf{p}) \|\mathbf{v}\|_1^2 \end{aligned}$$

Now, the assumptions introduced above are satisfied.

Next, we sketch the analysis of discretized problem. To discretize the equation choose a finite-dimensional space $V_h \subset V$ and a basis ϕ_1, \dots, ϕ_m of V_h . Usually, m is of order $1/h$. For the solution of the discretized equation one has to solve a linear system of equations which is uniquely solvable to the coercivity. To exploit Cea's Lemma we have to estimate $\text{dist}(\mathbf{u}, V_h)$. To do this one may choose the interpolation operator I_h defined by

$$I_h v := \sum_{i=1}^m v(\xi_i) \phi_i$$

where $0 = \xi_1 < \xi_2 < \dots < \xi_m = 1$ is a partition of $[0, 1]$. Then

$$\text{dist}(\mathbf{u}, V_h) \leq \|\mathbf{u} - I_h \mathbf{u}\|.$$

To estimate $\|\mathbf{u} - I_h \mathbf{u}\|$ is now a standard problem in the analysis of functions. □

5.4 Feasibility problems

Feasibility problems are ubiquitous in mathematical optimization. Given a constraint set, a solution to the feasibility problem is any point contained in the set. Frequently the constraint set can be realized as the intersection of a family of sets, each describing a simpler constraint. Thus, the convex feasibility problem can be formulated in a very simple way: Let C_1, \dots, C_m be a family of nonempty closed convex sets in a Hilbert space \mathcal{H} . The **convex feasibility problem** is:

$$\text{Given } C_1, \dots, C_m \text{ find a point } x \in C.$$

Historically, this problem is related to problem of guessing a starting point for solving a mathematical programming problem by the Simplex-algorithm.

To solve the feasibility problem one may try to find the projection $P_C(x^0)$ of a given „approximation“ $x^0 \in \mathcal{X}$. But it may be difficult to compute the metric projection P_C . Instead one may try to find the metric projection P_C by using the metric projections C_1, \dots, C_m which may be easier to compute in a cyclic way; see (??). That this idea may work might be seen in the simple example when C_1, C_2 are linear subspaces in \mathbb{R}^2 .

Linear systems

Consider the problem of solving the linear system

$$A\mathbf{x} = \mathbf{b} \tag{5.10}$$

where $A \in \mathbb{R}^{m,n}$ and $\mathbf{b} \in \mathbb{R}^m$. Let \mathbf{a}^i be the i th row of A and define

$$H_i := \{x \in \mathbb{R}^n : \langle \mathbf{a}^i | x \rangle = \mathbf{b}_i\}, \quad i = 1, \dots, m.$$

We assume $\mathbf{a}^i \neq \theta, i = 1, \dots, m$. Then each H_i is a hyperplane. Now, $A\mathbf{x} = \mathbf{b}$ if and only if $x \in H := \bigcap_{i=1}^m H_i$. The metric projection onto H_i is given by

$$P_{H_i}(x) := x + \frac{\mathbf{b}_i - \langle \mathbf{a}^i | x \rangle}{\|\mathbf{a}^i\|^2} \mathbf{a}^i \tag{5.11}$$

For numerical purposes, the computation of this projection is very pleasant since its computation requires vector arithmetic only.

The system (5.10) is called **consistent** if there exists a $z \in \mathbb{R}^n$ with $Az = \mathbf{b}$. In practice, the system (5.10) may be **inconsistent** due to the fact that it is overdetermined ($\mathbf{m} \gg \mathbf{n}$). As we will see in the next chapter, both cases can be handled by the cyclic projection method.

Suppose that instead we want to solve a linear system of inequalities

$$Ax \leq \mathbf{b} \quad (5.12)$$

where \leq is used componentwise. Then $Ax \leq \mathbf{b}$ if and only if $x \in \bigcap_{i=1}^m \hat{H}_i$ where \hat{H}_i are the half-spaces given by

$$\hat{H}_i := \{x \in \mathbb{R}^n : \langle \mathbf{a}^i | x \rangle \leq \mathbf{b}_i\}, \quad i = 1, \dots, m.$$

The metric projection onto \hat{H}_i is given by

$$P_{\hat{H}_i}(x) := \begin{cases} x & \text{if } \langle \mathbf{a}^i | x \rangle \leq \mathbf{b}_i \\ P_{H_i}(x) & \text{otherwise} \end{cases} \quad (5.13)$$

Again, there are many variations. For instance, we may consider the constrained linear system

$$Ax = \mathbf{b}, \quad x \geq \theta.$$

Here, in addition to \mathbf{m} affine constraints a cone-type constraint has to be met. The metric projection onto the cone $\mathbb{R}_+^n := \{x \in \mathbb{R}^n : x \geq \theta\}$ is a special case of the projection onto half-spaces.

Determining separating hyperplanes

Let C, D closed convex subsets in the Hilbert space \mathcal{H} such that $C \cap D = \emptyset$. A hyperplane $H = H_{\mathbf{y}, \mathbf{a}} := \{x \in \mathcal{H} : \langle \mathbf{y} | x \rangle = \mathbf{a}\}$ is called a **separating hyperplane** of C, D if

$$\langle \mathbf{y} | \mathbf{u} \rangle \leq \mathbf{a} \leq \langle \mathbf{y} | \mathbf{v} \rangle \text{ for all } \mathbf{u} \in C, \mathbf{v} \in D.$$

As we will see in the next chapter, the method of alternating projection may be used to construct two sequences $(\mathbf{u}^n)_{n \in \mathbb{N}}$ and $(\mathbf{v}^n)_{n \in \mathbb{N}}$ in C and D respectively and points $\mathbf{u} \in C, \mathbf{v} \in D$ such that

$$\mathbf{u} = \lim_n \mathbf{u}^n, \mathbf{v} = \lim_n \mathbf{v}^n \text{ and } \|\mathbf{u} - \mathbf{v}\| = \text{dist}(C, D).$$

We define $\mathbf{y} := \mathbf{u} - \mathbf{v}$, and $\mathbf{a} := t\mathbf{u} + (1-t)\mathbf{v}$ for some $t \in [0, 1]$ and define the set

$$H := \{x \in \mathcal{H} : \langle \mathbf{y} | x \rangle = \mathbf{a}\}.$$

Assumption: $\text{dist}(C, D) > 0$.

Then we must have $\mathbf{y} \neq \theta$ and H is a hyperplane. Further, $0 < \|\mathbf{y}\|^2 = \langle \mathbf{y} | \mathbf{u} \rangle - \langle \mathbf{y} | \mathbf{v} \rangle$ which implies $\langle \mathbf{y} | \mathbf{u} \rangle > \langle \mathbf{y} | \mathbf{v} \rangle$. As we know from Kolmogorov's criterion we have

$$\langle \mathbf{v} - \mathbf{u} | \mathbf{u}' - \mathbf{u} \rangle \leq 0, \quad \langle \mathbf{u} - \mathbf{v} | \mathbf{v}' - \mathbf{v} \rangle \leq 0 \text{ for all } \mathbf{u}' \in C, \mathbf{v}' \in D.$$

This implies

$$\langle \mathbf{y} | \mathbf{u}' \rangle \leq \langle \mathbf{y} | \mathbf{v} \rangle < \langle \mathbf{y} | \mathbf{u} \rangle \leq \langle \mathbf{y} | \mathbf{u}' \rangle \text{ for all } \mathbf{u}' \in C, \mathbf{v}' \in D,$$

and we conclude that H is a separating hyperplane.

Linear programming - a first view

Consider the **linear program (P)** given by

$$\text{Maximize } \langle \mathbf{c} | \mathbf{x} \rangle \text{ subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \quad (5.14)$$

where $\mathbf{A} \in \mathbb{R}^{m,n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$. Then (P) has a (symmetric) **dual program (D)** given by

$$\text{Minimize } \langle \mathbf{b} | \mathbf{y} \rangle \text{ subject to } \mathbf{A}^t \mathbf{y} \geq \mathbf{c}, \mathbf{y} \geq \mathbf{0}; \quad (5.15)$$

see for instance [7]. If the optimal value of either (P) or (D) is finite, then the optimal values of both (P) and (D) are finite and strong duality holds. That is, $\langle \mathbf{c} | \mathbf{x}^* \rangle = \langle \mathbf{b} | \mathbf{y}^* \rangle$ if and only if \mathbf{x}^* is optimal for (P) and \mathbf{y}^* is optimal for (D).

We define the constraints

$$\begin{aligned} \mathbf{C}_1 &:= \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^m : \langle \mathbf{c} | \mathbf{x} \rangle = \langle \mathbf{b} | \mathbf{y} \rangle\} && \text{(optimality)} \\ \mathbf{C}_2 &:= \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^m : \mathbf{A}\mathbf{x} \leq \mathbf{b}\} && \text{(primal feasibility)} \\ \mathbf{C}_3 &:= \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^m : \mathbf{x} \geq \mathbf{0}\} && \text{(primal feasibility)} \\ \mathbf{C}_4 &:= \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^m : \mathbf{A}^t \mathbf{y} \geq \mathbf{c}\} && \text{(dual feasibility)} \\ \mathbf{C}_5 &:= \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^m : \mathbf{y} \geq \mathbf{0}\} && \text{(dual feasibility)} \end{aligned}$$

Then $(\mathbf{x}^*, \mathbf{y}^*) \in \cap_{i=1}^5 \mathbf{C}_i$ if and only if \mathbf{x}^* is optimal for (P) and \mathbf{y}^* is optimal for (D). Since the variables \mathbf{x} and \mathbf{y} are related only by \mathbf{C}_1 it is convenient to consider the following three constraints: $\mathbf{C}_1, \mathbf{C}_2 \cap \mathbf{C}_4, \mathbf{C}_3 \cap \mathbf{C}_5$. \mathbf{C}_1 is a hyperplane, $\mathbf{C}_2 \cap \mathbf{C}_4$ is an intersection of halfplanes, $\mathbf{C}_3 \cap \mathbf{C}_5$ is the cone $\mathbb{R}_+^n \times \mathbb{R}_+^m$.

Semidefinite optimization

Consider the space $\mathbb{R}^{n,n}$ of square matrices. $\mathbb{R}^{n,n}$ can be endowed with the inner product

$$\langle \mathbf{A} | \mathbf{B} \rangle := \mathbf{A} \bullet \mathbf{B} := \sum_{i,j=1}^n \mathbf{a}_{ij} \mathbf{b}_{ij} = \text{trace}(\mathbf{B}^t \mathbf{A}), \mathbf{A} = (\mathbf{a}_{ij}), \mathbf{B} = (\mathbf{b}_{ij}).$$

Since $\mathbb{R}^{n,n}$ is finite-dimensional ($\dim \mathbb{R}^{n,n} = n^2$) $(\mathbb{R}^{n,n}, \langle \cdot | \cdot \rangle)$ is a Hilbert space. The inner product induces the so called **Frobenius norm**

$$\|\mathbf{A}\|_F := \langle \mathbf{A} | \mathbf{A} \rangle^{\frac{1}{2}}, \mathbf{A} \in \mathbb{R}^{n,n}.$$

The subset of symmetric matrices in $\mathbb{R}^{n,n}$ will be denoted by \mathcal{S}^n . The subset of positive semidefinite symmetric matrices is \mathcal{S}_+^n and \mathcal{S}_{++}^n is the set of symmetric positive definite matrices. Clearly, \mathcal{S}_+^n is a nonempty closed convex cone. To abbreviate the notions, we write occasionally $\mathbf{A} \succeq \mathbf{0}$ for $\mathbf{A} \in \mathcal{S}_+^n$.

Consider the primal semidefinite program:

$$\text{Minimize } \langle \mathbf{C} | \mathbf{X} \rangle \text{ subject to } \langle \mathbf{A}_i | \mathbf{X} \rangle = \mathbf{b}_i, i = 1, \dots, m, \mathbf{X} \succeq \mathbf{0}, \quad (5.16)$$

where $\mathbf{A} \in \mathbb{R}^{m,n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$.

The feasibility problem of this program may be reformulated as follows:

$$\text{Find } X \in \mathcal{S}_+^n \text{ with } \langle A_i | X \rangle = b_i, i = 1, \dots, m.$$

We set

$$C := \mathcal{S}_+^n, D := \{X \in \mathcal{S}^n : \langle A_i | X \rangle = b_i, i = 1, \dots, m\}.$$

Then we have to find a point in $C \cap D$.

The projection of a matrix X onto $C := \mathcal{S}_+^n$ can be found from the eigenvalue decomposition $X = \sum_{i=1}^n \lambda_i q_i q_i^t$ of X :

$$P_{\mathcal{S}_+^n}(X) = \sum_{i=1}^n \max\{0, \lambda_i\} q_i q_i^t$$

The projection of X onto on the affine set D is also easy to work out:

$$P_{A_i \bullet X = b_i, i=1, \dots, m}(X) = X - \sum_{i=1}^m u_i A_i$$

where u_i are found from the normal equations

$$Gu = (\langle A_i | X \rangle = b_i)_{i=1, \dots, m}, G = (G_{ij}), G_{ij} = \langle A_i | A_j \rangle, i, j = 1, \dots, m.$$

Remark 5.8. *As in the case of linear programming, we have a dual problem of semidefinite programs. In contrast to the linear programming, if the optimal value of either the primal or the dual program is finite, then the optimal values of both problems are not forced to be equal; a duality gap may exist. But there are well known conditions which are sufficient for the non-existence of a duality gap. If there is no duality gap the feasibility problem for the primal and dual problem together with the optimality condition may be considered as in the linear programming context. ■*

An interesting application for the semidefinite approach is the the so called **educational testing problem**; see [16, 1]. Such optimization problems arise in statistics.

Given a matrix $A \in \mathcal{S}_+^n$, find out how much can be subtracted from the diagonal of A and still retain a matrix in \mathcal{S}_+^n .

This can be expressed as

$$\text{Maximize } \langle e | d \rangle \text{ subject to } A - \text{Diag}(d) \in \mathcal{S}_+^n, d \geq 0, \quad (5.17)$$

where $e := (1, 1, \dots, 1) \in \mathbb{R}^n$.

Fletcher [23, 24] describes methods to solve the problem by „classical“ methods of optimization. W. Glunt [25] applies successfully the idea described in the following subsection.

Linear programming - a second view

We consider a variation of the linear programming problem (5.16).

$$\text{Minimize } \langle \mathbf{e} | \mathbf{x} \rangle \text{ subject to } \mathbf{x} \in \bigcap_{i=1}^m C_i, \quad (5.18)$$

where $\mathbf{e} = (1, 1, \dots, 1) \in \mathbb{R}^n$ and C_1, \dots, C_m are closed convex subsets of \mathbb{R}^n . Clearly, the feasible set $C := \bigcap_{i=1}^m C_i$ is a generalized constraint set; compare with (5.16). We call the problem (5.18) a **linear convex programming problem**. To relate this problem to alternate projection methods following an idea of W. Glunt [25] we introduce the hyperplane

$$L_\tau := \{\mathbf{y} \in \mathbb{R}^n : \langle \mathbf{e} | \mathbf{y} \rangle = \tau\}.$$

If τ is chosen such that

$$\tau < \min_{\mathbf{u} \in C} \langle \mathbf{e} | \mathbf{u} \rangle$$

then the sets C and L_τ are disjoint. Given a $\mathbf{f} \in \mathbb{R}^n$ one may consider the problem

$$\text{Minimize } \|\mathbf{f} - \mathbf{x}\| \text{ subject to } \mathbf{x} \in C \cap L_\tau. \quad (5.19)$$

which has no feasible point. Nevertheless, one may consider the alternate projection methods to try to compute $P_{C \cap L_\tau}(\mathbf{f})$. What we get by this procedure are two sequences $(\mathbf{u}^n)_{n \in \mathbb{N}}, (\mathbf{y}^n)_{n \in \mathbb{N}}$, belonging to C and L_τ respectively. As we show below following [15], these sequences converge to points $\mathbf{u} \in C$ and $\mathbf{y} \in L_\tau$ such that $\|\mathbf{u} - \mathbf{y}\|$ attains the minimum distance between C and L_τ . It can be deduced from the relationship of L_τ and $\langle \mathbf{e} | \cdot \rangle$ that \mathbf{u} solves the problem (5.18).

Theorem 5.9. *Let \mathcal{H} be a Hilbert space and let C, D be nonempty closed convex subsets of \mathcal{H} . Then under the assumption*

$$\text{one of the sets } C, D \text{ is compact} \quad (5.20)$$

there exist $\mathbf{u} \in C$ and $\mathbf{v} \in D$ respectively such that

$$\|\mathbf{u} - \mathbf{w}\| = \inf_{\mathbf{u}' \in C, \mathbf{y}' \in D} \|\mathbf{u}' - \mathbf{y}'\|$$

Proof:

Let us assume that C is compact.

Obviously, there exist a sequence $(\mathbf{u}^n)_{n \in \mathbb{N}}$ in C and a sequence $(\mathbf{v}^n)_{n \in \mathbb{N}}$ in D with

$$\|\mathbf{u}^n - \mathbf{v}^n\| - \frac{1}{n} \leq \mathbf{d} := \inf_{\mathbf{u} \in C, \mathbf{v} \in D} \|\mathbf{u} - \mathbf{v}\|, n \in \mathbb{N}.$$

Since C is compact $(\mathbf{u}^n)_{n \in \mathbb{N}}$ contains a convergent subsequence $(\mathbf{u}^{n_k})_{k \in \mathbb{N}}$; let $\mathbf{u} := \lim_k \mathbf{u}^{n_k}$. Since C is compact and therefore closed, $\mathbf{u} \in C$. Then $(\mathbf{v}^{n_k})_{k \in \mathbb{N}}$ is bounded and contains therefore a weakly convergent subsequence $(\mathbf{v}^{n_{k_l}})_{l \in \mathbb{N}}$; let $\mathbf{v} = \mathbf{w} - \lim_l \mathbf{v}^{n_{k_l}}$. Since D is closed and convex, D is weakly closed and therefore $\mathbf{v} \in D$. Now $(\mathbf{u}^{n_{k_l}} - \mathbf{v}^{n_{k_l}})_{l \in \mathbb{N}}$ converges weakly to $\mathbf{u} - \mathbf{v}$. Since the norm in \mathcal{H} is weakly lower semicontinuous, $\mathbf{d} = \|\mathbf{u} - \mathbf{v}\|$. ■

Cheney and Goldstein have shown in 1959 [15] that under the assumption of Theorem 5.9 the alternate projection method produces a vector \mathbf{w} which satisfies

$$\mathbf{w} = \mathbf{u} - \mathbf{v}, \mathbf{u} \in C, \mathbf{v} \in D, \|\mathbf{w}\| = \text{dist}(C, D).$$

5.5 Approximation of correlation matrices

A correlation matrix is a symmetric positive semidefinite matrix with unit diagonal. Correlation matrices occur in several areas of numerical linear algebra, including preconditioning of linear systems and error analysis of Jacobi methods for the symmetric eigenvalue problem. The term *correlation matrix* comes from statistics, since a matrix whose (i, j) entry is the correlation coefficient between two random variables x_i and x_j is a symmetric positive semidefinite matrix with unit diagonal. In stock research sample correlation matrices constructed from vectors of stock returns are used for predictive purposes. Unfortunately, on any day when an observation is made data is rarely available for all the stocks of interest. One way to deal with this problem is to compute the sample correlations of pairs of stocks using data drawn only from the days on which both stocks have data available. The resulting matrix of correlations will be only an approximate correlation matrix, because it has been built from inconsistent data sets.

The problem we consider is for a matrix $A \in \mathcal{S} := \mathcal{C}_1 \cap \mathcal{S}_+^n$ to compute the best approximation of A in \mathcal{S} with respect to the Frobenius-norm:

$$\text{Minimize } \|X - A\|_F \text{ subject to } A \in \mathcal{C}_1 \cap \mathcal{S}_+^n$$

This leads to the computation $P_{\mathcal{S}}(A)$. Since \mathcal{C}_1 and \mathcal{S}_+^n are both closed convex subsets of $\mathbb{R}^{n,n}$, so is their intersection. Therefore the metric projection $P_{\mathcal{S}}$ exists. Moreover, each $P_{\mathcal{S}}(A)$ is a singleton. The alternate projection method may be used to compute $P_{\mathcal{S}}$; see [29, 30, ?, 50]

In practice, it is reasonable to fit the norm in $\mathbb{R}^{n,n}$ to some observations concerning the quality of the correlation data. One can do this by choosing a matrix $W \in \mathcal{S}_{++}^n$. This W induces the inner product $\langle A|B \rangle_W := \langle A|WB \rangle_F$ and the norm $\|A\|_W := \|W^{\frac{1}{2}}A\|_F$.

5.6 Extremal property of splines

Here we consider the interpolation of „difficult to evaluate functions“ by „simpler“ functions. We follow here mainly [4, 9]

The curvature of a curve given by the parametrization

$$\gamma : [a, b] \ni t \longmapsto (t, f(t)) \in \mathbb{R}^2$$

with a twice continuous differentiable function f is given by

$$\kappa_{\gamma}(t) := \frac{f''(t)}{(1 + f'(t)^2)^{\frac{3}{2}}}, \quad t \in [a, b]. \quad (5.21)$$

From a physical point of view, the integral

$$\int_a^b \frac{f''(t)^2}{(1 + f'(t)^2)^3} dt \quad (5.22)$$

is the strain energy of a flexible thin beam or draftman's *spline*¹ forced to pass through prescribed points. Consider a partition $\tau : \mathbf{a} = \mathbf{t}_0 < \mathbf{t}_1 < \dots < \mathbf{t}_m = \mathbf{b}$. We want to find a curve of least strain energy which interpolates the given data $\mathbf{y}_0, \dots, \mathbf{y}_m \in \mathbb{R}$. The criterion in (5.22) is a rather nonlinear one. Therefore it is appropriate to approximate the strain energy. Historically and from the practical point of view a very successful approximation is the following:

$$\int_a^b f''(t)^2 dt \quad (5.23)$$

If we summarize the interpolation conditions by

$$W(\mathbf{y}) := \{f \in C^2[\mathbf{a}, \mathbf{b}] : f(\mathbf{t}_i) = \mathbf{y}_i, i = 0, \dots, m\}$$

we may state the interpolation problem by

$$\text{Minimize } \int_a^b f''(t)^2 dt \text{ subject to } f \in W(\mathbf{y}) \quad (5.24)$$

Unfortunately, $C^2[\mathbf{a}, \mathbf{b}]$ is not complete in the norm

$$\|f\|_c := \left(\int_a^b |f(t)|^2 + |f''(t)|^2 dt \right)^{\frac{1}{2}}.$$

Therefore we consider a completion $H^2[\mathbf{a}, \mathbf{b}]$ of $C^2[\mathbf{a}, \mathbf{b}]$ defined as follows:

$$AC[\mathbf{a}, \mathbf{b}] := \{f \in C[\mathbf{a}, \mathbf{b}] : f(t) = \alpha + \int_a^t g(s) ds, t \in [\mathbf{a}, \mathbf{b}], \text{ with } g \in L_1[\mathbf{a}, \mathbf{b}], \alpha \in \mathbb{R}\},$$

$$H^2[\mathbf{a}, \mathbf{b}] := \{f \in C^1[\mathbf{a}, \mathbf{b}] : f'(t) \in AC[\mathbf{a}, \mathbf{b}], f' = \alpha + \int_a^t g(s) ds, t \in [\mathbf{a}, \mathbf{b}], g \in L^2[\mathbf{a}, \mathbf{b}]\}.$$

The function g in the definition of $H^2[\mathbf{a}, \mathbf{b}]$ above is called the second derivative of f in a generalized sense and we write $f'' := g$. This space $H^2[\mathbf{a}, \mathbf{b}]$ is a special case of the family of **Sobolev spaces**. It is a Hilbert space under the inner product

$$\langle f, h \rangle_{2,2} := \int_a^b (f(t)h(t) + f''(t)h''(t)) dt.$$

A half norm in $H^2[\mathbf{a}, \mathbf{b}]$ is given by

$$\nu(f)_{2,2} := \left(\int_a^b |f''(t)|^2 dt \right)^{\frac{1}{2}}.$$

Now we reformulate the minimization problem above as follows:

$$\text{Minimize } \nu(f)_{2,2} \text{ subject to } f \in V(\mathbf{y}) \quad (5.25)$$

where we have adapted the interpolating conditions to the new context:

$$V(\mathbf{y}) := \{f \in H^2[\mathbf{a}, \mathbf{b}] : f(\mathbf{t}_i) = \mathbf{y}_i, i = 0, \dots, m\}$$

¹The term „spline“ was introduced by Schoenberg 1946.

The problem (5.24) has now a formulation which is appropriate to analyze the problem by methods prepared in the last chapter. The problem that $\nu_{2,2}$ is not a norm is not very serious because in general the interpolating conditions enforces uniqueness of a solution of (5.25).

One can show that such a solution f is a C^2 -function and on each interval $[t_i, t_{i+1}]$, $i = 0, \dots, m-1$, equal to a polynomial of degree at least three. Therefore, a solution of the minimization problem (5.25) is called a **interpolating cubic spline** and the fact that such a spline is a solution of (5.25) is called the **extremal property** of the spline. On the boundary of the interval $[a, b]$ we observe the **natural boundary conditions** $f''(a) = f''(b) = 0$. If we count the degrees of freedom which such a function contains then we obtain the number $4m$. On the other hand we have to obey $3(m-1) + (m+1)$ requirements (continuity, interpolation conditions). There are two degrees of freedom left which generate the natural boundary conditions. These natural boundary conditions may be replaced by two other conditions; see below.

Let us consider a more general problem formulation. We choose Hilbert spaces \mathcal{H}, \mathcal{K} , a linear continuous operator $R : \mathcal{K} \rightarrow \mathcal{H}$, a convex subset M of \mathcal{K} , and an affine closed subspace V of \mathcal{K} . Then we want to solve

$$\text{Minimize } \|R(x)\|_{\mathcal{H}} \text{ subject to } x \in V \cap M \quad (5.26)$$

This problem may be considered as an approximation problem:

$$\text{Find in } \mathcal{H} \text{ the best approximation } w \text{ of } \theta \text{ in } C := R(V \cap M). \quad (5.27)$$

The set C is convex but the existence of a solution is in doubt since it is not clear whether C is a closed. We introduce now a description of V with the intention to analyze this question.

Let \mathcal{Y} be a Banach space and let $A : \mathcal{K} \rightarrow \mathcal{Y}$ be a linear continuous mapping with nullspace $\ker(A)$. Moreover, we assume that A describes the set V as follows: $V = A^{-1}(y)$ with a given $y \in \mathcal{Y}$. On this way we describe interpolating conditions: $x \in \mathcal{K}$ is „interpolating the data y “ if $Ax = y$. With the set M we describe constraints of non-interpolatory character; for example: $M = \{x \in \mathcal{K} : x \geq \theta\}$ where \geq is a (partial) order in \mathcal{K} .

Lemma 5.10. *Let $R : \mathcal{K} \rightarrow \mathcal{H}$ be a surjective linear operator, let $A : \mathcal{K} \rightarrow \mathcal{Y}$ be a linear operator. If $\ker(R) \cap \ker(A) = \{\theta\}$ and if $\ker(R)$ is finite-dimensional then*

$$\|\cdot\|_{(R,A)} : \mathcal{K} \ni x \mapsto (\|R(x)\|^2 + \|A(x)\|^2)^{\frac{1}{2}} \in \mathbb{R} \quad (5.28)$$

defines a norm in \mathcal{K} which is equivalent to the norm in \mathcal{K} .

Proof:

Obviously, the expression in (5.28) is a norm due to the assumption $\ker(R) \cap \ker(A) = \{\theta\}$. We have

$$\|R(x)\|^2 + \|A(x)\|^2 \leq \|R\|^2 \|x\|^2 + \|A\|^2 \|x\|^2$$

which shows that there exists a constant $c_1 > 0$ with

$$\|x\|_{(R,A)} \leq c_1 \|x\| \text{ for all } x \in \mathcal{K}.$$

Consider the mapping $L : \mathcal{K} \ni x \mapsto (Rx, Ax) \in \mathcal{H} \times \mathcal{Y}$. Clearly, $\mathcal{H} \times \mathcal{Y}$ is a Banach space and L is a linear continuous operator. We show that the range $\text{ran}(L)$ is closed.

Let $(L(x^n))_{n \in \mathbb{N}}$ be a sequence in $\text{ran}(L)$ with $\lim_n R(x^n) = w, \lim_n A(x^n) = z$. We decompose each x^n as $x^n = u^n + v^n \in \ker(R) + \ker(R)^\perp$. We have $\lim_n R(v^n) = w$. Since $\hat{R} := R|_{\ker(R)^\perp}$ is linear, continuous and bijective we conclude that $(v^n)_{n \in \mathbb{N}}$ is a Cauchy sequence and therefore convergent. Let $v := \lim_n v^n \in \ker(R)^\perp$.

Assume that $(x^n)_{n \in \mathbb{N}}$ is not bounded. Then we conclude from

$$\frac{x^n}{\|x^n\|} = \frac{u^n}{\|x^n\|} + \frac{v^n}{\|x^n\|}$$

that $(\frac{u^n}{\|x^n\|})_{n \in \mathbb{N}}$ is a bounded sequence in $\ker(R)$ and a convergent subsequence $(\frac{u^{n_k}}{\|x^{n_k}\|})_{k \in \mathbb{N}}$ exists. Clearly, $u := \lim_k \frac{u^{n_k}}{\|x^{n_k}\|} \in \ker(R)$. Since $(Ax^{n_k})_{k \in \mathbb{N}}, (Av^{n_k})_{k \in \mathbb{N}}$ are bounded we obtain $Au = \theta$, which implies $u \in \ker(R) \cap \ker(A)$ and finally $u = \theta$. But this is contradicting the fact that $\frac{x^{n_k}}{\|x^{n_k}\|} = 1, k \in \mathbb{N}$.

Now we know that $(x^n)_{n \in \mathbb{N}}$ is bounded which implies that $(u^n)_{n \in \mathbb{N}}$ is bounded too. Hence, $(u^n)_{n \in \mathbb{N}}$ has a convergent subsequence and since $(v^n)_{n \in \mathbb{N}}$ is convergent, $(x^n)_{n \in \mathbb{N}}$ has a convergent subsequence too. Each cluster point x of $(x^n)_{n \in \mathbb{N}}$ leads to $(w, z) = (Rx, Ax)$. Now, the range of L is closed. Then L is a linear, continuous bijective mapping from \mathcal{H} onto $\text{ran}(L)$. Therefore $L^{-1} : \text{ran}(L) \rightarrow \mathcal{K}$ exists and is continuous. From

$$\|x\| = \|L^{-1}(Rx, Ax)\| \leq \|L^{-1}\| \|(Rx, Ax)\| = \|L^{-1}\| (\|Rx\|^2 + \|Ax\|^2)^{\frac{1}{2}} = \|L^{-1}\| \|x\|_{(R,A)}$$

we conclude

$$c_1^{-1} \|x\|_{(R,A)} \leq \|x\| \leq \|L^{-1}\| \|x\|_{(R,A)} \text{ for all } x \in \mathcal{K}. \quad \blacksquare$$

Theorem 5.11. *Let $R : \mathcal{K} \rightarrow \mathcal{H}$ be a surjective linear operator, let $A : \mathcal{K} \rightarrow \mathcal{Y}$ be a linear operator and let $y \in \mathcal{Y}$. We set $V := A^{-1}(y)$ and let M be a closed convex subset in \mathcal{K} . If $\ker(R) \cap \ker(A) = \{\theta\}$ and if $\ker(R)$ is finite-dimensional then $R(V \cap M)$ is closed and the spline-extremal problem*

$$\text{Minimize } \|Rx\| \text{ subject to } x \in V \cap M. \quad (5.29)$$

has a uniquely determined solution $w \in V \cap M$.

Proof:

Set $S := R(V \cap M)$. Obviously, S is convex. We want to show that S is closed. Let $(x^n)_{n \in \mathbb{N}}$ a sequence in $V \cap M$ with $\lim_n R(x^n) = u$. Then $\lim_n (Rx^n, Ax^n) = (u, y)$. Therefore $(Rx^n, Ax^n)_{n \in \mathbb{N}}$ is a Cauchy sequence in $\mathcal{H} \times \mathcal{Y}$ with respect to the norm $\|\cdot\|_{(R,A)}$; see Lemma 5.10. Again, by Lemma 5.10 $(x^n)_{n \in \mathbb{N}}$ is a Cauchy sequence in \mathcal{K} and therefore convergent. Let $x = \lim_n x^n$. Then $x \in M, Ax = y, Rx = u$.

Now, there exists a uniquely determined $v \in S$ with

$$\|v\| = \inf_{v' \in S} \|v'\|$$

We have $v = Rw, w \in V \cap M$. Suppose that there is another $w' \in V \cap M$ with $v = Rw'$. Then $w - w' \in \ker(R) \cap \ker(A)$ and by assumption, $w = w'$. \blacksquare

Next we want to characterize the solution w in Theorem 5.11. We know from Theorem 3.1 (Kolmogorov's criterion)

$$\langle \mathbf{R}w | \mathbf{R}v - \mathbf{R}w \rangle \geq 0 \text{ for all } v \in V \cap M. \quad (5.30)$$

This is a variational inequality in the space \mathcal{H} . Using the adjoint \mathbf{R}^* of \mathbf{R} we obtain a variational inequality in \mathcal{K} :

$$\langle \mathbf{R}^* \mathbf{R}w | v - w \rangle \geq 0 \text{ for all } v \in V \cap M. \quad (5.31)$$

The ideas in the section on variational inequalities can be applied to (5.32).

If we have no constraints M , i.e. if $M = \mathcal{K}$, then the variational inequality (5.32) becomes

$$\mathbf{R}w \in \mathbf{R}(\ker(\mathbf{A}))^\perp. \quad (5.32)$$

Lemma 5.12. *Let $\mathbf{R} : \mathcal{K} \rightarrow \mathcal{H}$ be a surjective linear operator, let $\mathbf{A} : \mathcal{K} \rightarrow \mathcal{Y}$ be a linear surjective operator and let $\mathbf{y} = \mathbf{A}x^0 \in \mathcal{Y}$; we set $V(\mathbf{y}) := x^0 + \ker(\mathbf{A})$. Let $\mathbf{R}(w) \in \mathbf{R}(\ker(\mathbf{A}))^\perp$ with $w \in V(\mathbf{y})$. We assume*

$$\dim(\ker(\mathbf{R})) < \infty, \dim \mathcal{Y} < \infty, \ker(\mathbf{R}) \cap \ker(\mathbf{A}) = \{\emptyset\}. \quad (5.33)$$

Then

- (a) $\mathbf{R}^*(\mathbf{R}(\ker(\mathbf{A})))^\perp = \ker(\mathbf{A}) \cap \ker(\mathbf{R})$.
- (b) $m + 1 = \dim(\mathbf{R}(\ker(\mathbf{A})))^\perp = \dim(\ker(\mathbf{A}) \cap \ker(\mathbf{R}))$.
- (c) $\ker(\mathbf{A})^\perp$ is finite-dimensional.
- (d) $\dim(\ker(\mathbf{A})^\perp) = \dim(\ker(\mathbf{R})) + \dim(\ker(\mathbf{A})^\perp \cap \ker(\mathbf{R})^\perp)$.
- (e) If $\mathbf{b}_1, \dots, \mathbf{b}_m$ is a basis of $\mathbf{R}(\ker(\mathbf{A}))^\perp$ then

$$\langle \mathbf{R}w | \mathbf{b}_i \rangle = \langle \mathbf{y} | \eta_i \rangle, \quad i = 1, \dots, m, \quad (5.34)$$

where $\mathbf{R}^* \mathbf{b}_i = \mathbf{A}^* \eta_i, i = 1, \dots, m$.

Proof:

Ad (a) This can easily be shown.

Ad (b) Follows from the fact that \mathbf{R}^* is injective.

Ad (c) Clearly, $\dim(\ker(\mathbf{A})^\perp) = \dim(\text{ran}(\mathbf{A})) < \infty$.

Ad (d) We have

$$\ker(\mathbf{A})^\perp = \ker(\mathbf{R}) \cap \ker(\mathbf{A})^\perp \oplus \ker(\mathbf{R})^\perp \cap \ker(\mathbf{A})^\perp = \ker(\mathbf{R}) \oplus \ker(\mathbf{R})^\perp \cap \ker(\mathbf{A})^\perp$$

since $\ker(\mathbf{R}) \subset \ker(\mathbf{A})^\perp$ due to the assumption $\ker(\mathbf{R}) \cap \ker(\mathbf{A}) = \{\emptyset\}$.

Ad (e) We have $\mathbf{R}^* \mathbf{b}_i \in \ker(\mathbf{R})^\perp \cap \ker(\mathbf{A})^\perp$. Hence there exist η_i with $\mathbf{R}^* \mathbf{b}_i = \mathbf{A}^* \eta_i, i = 1, \dots, m$. Now, $\mathbf{R}w = \mathbf{R}(x^0 + u)$ with $x^0 \in V(\mathbf{y}) \cap \ker(\mathbf{A})^\perp, u \in \ker(\mathbf{A})$. Then for $i = 1, \dots, m$

$$\langle \mathbf{R}w | \mathbf{b}_i \rangle = \langle \mathbf{R}x^0 | \mathbf{b}_i \rangle = \langle x^0 | \mathbf{R}^* \mathbf{b}_i \rangle = \langle x^0 | \mathbf{A}^* \eta_i \rangle = \langle \mathbf{A}x^0 | \eta_i \rangle = \langle \mathbf{y} | \eta_i \rangle.$$

■

We conclude from Lemma 5.12 that the solution \mathbf{w} of the spline-problem can be computed as follows: Solve the linear system

$$\mathbf{B}\mathbf{a} = \mathbf{r} \text{ with } \mathbf{B} := (\langle \mathbf{b}_i | \mathbf{b}_j \rangle) \in \mathbb{R}^{m,m}, \mathbf{r} = (\langle \mathbf{y} | \boldsymbol{\eta}_j \rangle) \in \mathbb{R}^m \quad (5.35)$$

with solution $\mathbf{a} \in \mathbb{R}^m$ and find \mathbf{w} as follows:

$$\mathbf{w} := \mathbf{R}^{-1} \mathbf{g} \in \ker(\mathbf{R})^\perp \text{ where } \mathbf{g} := \sum_{i=1}^m \mathbf{a}_i \mathbf{b}_i. \quad (5.36)$$

Notice that $\mathbf{w} \in \ker(\mathbf{R})^\perp$ in (5.36) is well-defined due to the assumption $\ker(\mathbf{R}) \cap \ker(\mathbf{A}) = \{\boldsymbol{\theta}\}$.

Let us come back to the interpolatory spline at the beginning of this section. Here

$$\begin{aligned} \mathcal{K} &:= \mathcal{H}^2[\mathbf{a}, \mathbf{b}], \mathcal{H} = \mathcal{L}_2[\mathbf{a}, \mathbf{b}], \mathbf{R}: \mathcal{K} \ni f \longmapsto f'' \in \mathcal{H}, \\ \mathcal{Y} &:= \mathbb{R}^{n+4}, \mathbf{A}: \mathcal{K} \ni f \longmapsto (f'(\mathbf{a}), f(\mathbf{a}), f(\mathbf{t}_1), \dots, f(\mathbf{t}_n), f(\mathbf{b}), f'(\mathbf{b})) \in \mathcal{Y}, \end{aligned}$$

where $\tau: \mathbf{a} = \mathbf{t}_0 < \mathbf{t}_1 < \dots < \mathbf{t}_n < \mathbf{t}_{n+1} = \mathbf{b}$ is a partition of $[\mathbf{a}, \mathbf{b}]$. Then all the requirements of the more general problem above are satisfied, especially the assumption $\ker(\mathbf{R}) \cap \ker(\mathbf{A}) = \{\boldsymbol{\theta}\}$. A basis $\mathbf{b}_0, \dots, \mathbf{b}_m$ in $\mathcal{R}(\ker(\mathbf{A}))^\perp$ is given by

$$\begin{aligned} \mathbf{b}_0(\mathbf{t}) &:= \begin{cases} \frac{\mathbf{t}_1 - \mathbf{t}}{\mathbf{t}_1 - \mathbf{t}_0} & , \mathbf{t} \in [\mathbf{t}_1, \mathbf{t}_0] \\ 0 & , \text{else} \end{cases} \\ \mathbf{b}_i(\mathbf{t}) &:= \begin{cases} \frac{\mathbf{t} - \mathbf{t}_i}{\mathbf{t}_{i-1} - \mathbf{t}_i} & , \mathbf{t} \in [\mathbf{t}_{i-1}, \mathbf{t}_i] \\ \frac{\mathbf{t}_{i+1} - \mathbf{t}}{\mathbf{t}_{i+1} - \mathbf{t}_i} & , \mathbf{t} \in [\mathbf{t}_i, \mathbf{t}_{i+1}], i = 1, \dots, n, \\ 0 & , \text{else} \end{cases} \\ \mathbf{b}_{n+1}(\mathbf{t}) &:= \begin{cases} \frac{\mathbf{t} - \mathbf{t}_n}{\mathbf{t}_{n+1} - \mathbf{t}_n} & , \mathbf{t} \in [\mathbf{t}_n, \mathbf{t}_{n+1}] \\ 0 & , \text{else} \end{cases} \end{aligned}$$

Since these basis-functions are piecewise linear functions the resulting solution of the spline-problem is piecewise a polynomial of degree three.

In order to find a nonlinear framework for the interpolatory splines one can replace the optimization criterion as follows:

$$\text{Minimize } \phi(\mathbf{R}\mathbf{x}) \text{ subject to } \mathbf{x} \in \mathbf{V} \cap \mathbf{M}, \quad (5.37)$$

where ϕ is a functional defined on \mathcal{H} . Since the optimization criterion has to be well-defined one may \mathcal{K}, \mathcal{H} adapt to this requirement. This can be done in the concrete case by replacing the space $\mathcal{L}_2[\mathbf{a}, \mathbf{b}]$ by an **Orlicz-type space**. In this framework the spline problem has been discussed in [4, 5, 9]. For example, one can characterize on this way piecewise rational functions by an extremal property.

Another approach is considered by B.J. McCartin [42]. Starting from the analogy of a cubic spline to a beam a tension term is added to the governing differential equation so giving rise to the exponential spline.

5.7 Variational inequalities

An obstacle problem

Let Ω be a bounded subset in \mathbb{R}^2 with boundary $\partial\Omega$. Consider the problem

$$-\Delta u = f, \text{ in } \Omega, \quad u = 0 \text{ in } \partial\Omega, \quad u \geq g \text{ in } \Omega \quad (5.38)$$

Here f, g are given functions in Ω and $\partial\Omega$, respectively. g is called the **obstacle**. We reformulate this problem by introducing the Sobolev space $V := H_0^1(\Omega)$, the operator $A : V \rightarrow V^*$, $\langle A(u), v \rangle := \int_{\Omega} \nabla u \nabla v dx$ and the functional $b \in V^*$, $\langle b, v \rangle := \int_{\Omega} f v dx$. The problem (5.38) becomes

$$\langle A(u) - b, v - u \rangle \leq 0 \text{ for all } v \in C, \quad (5.39)$$

where the obstacle is now described by

$$C := \{w \in V : w \geq g\}.$$

Solving variational inequalities by projection methods

Given a nonempty closed and convex subset C of a Hilbert space \mathcal{H} and a mapping $F : C \rightarrow \mathcal{H}$, the **variational inequality problem**² is to find a vector $x^* \in C$ such that

$$\langle F(x^*) | x - x^* \rangle \geq 0 \text{ for all } x \in C. \quad (5.40)$$

Variational inequalities provide a tool to formulate various equilibrium problems. Several well known problems, such as systems of nonlinear equations, first order conditions for linear and nonlinear optimization problems, and complementarity problems, are special cases of variational inequalities.

A geometric interpretation of a solution of a variational inequality is that x^* is a solution if and only if $F(x^*)$ makes a non-obtuse angle with all the feasible directions going into C from x^* . Kolmogorov's criterion is a specific case of this variational inequality. If a solution x^* belongs to the interior of C then obviously $F(x^*) = \theta$.

Let us just sketch the variational inequality which describes a necessary condition in optimization. Consider the optimization problem

$$\text{Minimize } f(x) \text{ subject } x \in C.$$

Here, f is a differentiable function from $\mathbb{R}^n \rightarrow \mathbb{R}$ and C is a convex subset of \mathbb{R}^n . As we know from analysis a necessary condition for a solution x^* of the minimization problem is the variational inequality

$$\langle \nabla f(x^*) | x - x^* \rangle \geq 0 \text{ for all } x \in C$$

Setting $F(x) := \nabla f(x)$ we obtain a variational inequality.

Theorem 5.13. *Let \mathcal{H} be a Hilbert space and let C be nonempty closed convex subset of \mathcal{H} . Then for $x^* \in C$ the following conditions are equivalent:*

²Variational inequalities were introduced by Hartman and Stampacchia in 1966 for the study of partial differential equations with applications in the field of mechanics.

(a) x^* is a solution of the variational inequality (5.40).

(b) x^* is a fixed point of $\mathcal{H} \ni x \mapsto P_C(x - \tau F(x))$ for each $\tau > 0$.

Proof:

(a) \implies (b) Let $\tau > 0$. Define $u := x^* - F(x^*)$. Then $\langle x^* - u | v - x^* \rangle \geq 0$ for all $v \in C$. Kolmogorov's criterion implies $x^* = P_C(u) = P_C(x^* - \tau F(x^*))$.

(b) \implies (a) We have $\langle x^* - (x^* - \tau F(x^*)) | v - x^* \rangle \geq 0$ for all $v \in C$. ■

Using Theorem 5.13 we may solve the variational inequality by solving the fixed point equation

$$x = P_C(x - \tau F(x)) \tag{5.41}$$

Since P_C is nonexpansive we can use the fixed point theory of nonexpansive mappings. If the constraint set C can be decomposed into a intersection of sets, $C = \bigcap_{i=1}^m C_i$, we may use a cyclic projection method to solve the fixed point equation.

Common solutions of variational inequalities

Let \mathcal{H} be a Hilbert space. Let there be given for each $i = 1, \dots, m$

- a mapping $A_i : \mathcal{H} \rightrightarrows \mathcal{H}$,
- a nonempty closed convex subset C_i .

Then one may consider the problem (CVI):

Find a point $x \in \bigcap_{i=1}^m C_i$ such that for each $i = 1, \dots, m$ there exists $u_i \in A_i(x)$ satisfying

$$\langle u_i | v - x \rangle \geq 0 \text{ for all } v \in C_i, i = 1, \dots, m. \tag{5.42}$$

5.8 A proof of Schauder's fixed point theorem

Theorem 5.14 (Fixed point Theorem of Schauder, 1930/Second Version). *Let X be a Banach space and let M be a nonempty compact convex subset of \mathcal{X} . If $F : M \rightarrow M$ is continuous then F possesses a fixed point.*

Proof:

As the closed linear span of the compact set M is a separable normed space, we may assume, without loss of generality, that \mathcal{X} is separable. Moreover, since a separable Banach space can be renormed such that the norm is equivalent to the given norm and strictly convex, we can assume that the norm in \mathcal{X} is strictly convex; see Koethe [34]. Notice that the existence of a fixed point is not influenced by renorming. For each $n \in \mathbb{N}$, let S_n be a finite subset of M such that $M \subset S_n + B_{1/n}$. Notice that the choice of S_n is possible since M is totally bounded³. Let C_n be the closed convex hull of S_n . Then C_n is a compact convex subset of M ; see the Lemma 5.15 below. Thus, P_{C_n} is single-valued and continuous. Let Z_n be the linear span of S_n . Then $\dim Z_n < \infty$. Define $F_n : Z_n \rightarrow Z_n$ by $F_n(x) := P_{C_n}(x)$. By Brouwer's fixed point theorem F_n has a fixed point $x^n \in C_n$. Now

$$\|F_n(x) - F(x)\| \leq 1/n \text{ for all } x \in M.$$

³A set in a metric space is totally bounded if it can be covered by finitely many balls of any fixed size.

Therefore, $\|F(x^n) - x^n\| \leq 1/n$, $n \in \mathbb{N}$. Now, it is clear that every cluster point of $(x^n)_{n \in \mathbb{N}}$, and such cluster points exist since M is compact, is a fixed point of F . ■

Lemma 5.15 (Mazur, 1934). *Let X be a normed space and let $M \subset X$ be such that \overline{M} is compact. Then $\overline{\text{co}}(M)$ is compact too.*

Proof:

Let $\varepsilon > 0$. Since \overline{M} is compact there exist x^1, \dots, x^n in M with

$$M \subset \bigcup_{i=1}^n B_{\varepsilon/2}(x^i).$$

Let $x \in \text{co}(M)$. Then this x has a presentation

$$x = \sum_{i=1}^m a_i y^i, \quad y^1, \dots, y^m \in M, \quad a_1, \dots, a_m \geq 0, \quad \sum_{i=1}^m a_i = 1.$$

Now, we have m indices j_1, \dots, j_m such that $y^i - x^{j_i} \in B_{\varepsilon/2}$, $i = 1, \dots, m$. Therefore we obtain

$$x - \sum_{i=1}^m a_i x^{j_i} = \sum_{i=1}^m a_i (y^i - x^{j_i}) \in B_{\varepsilon/2}.$$

This implies $K := \text{co}(M) \subset \text{co}(\{x^1, \dots, x^n\}) + B_{\varepsilon/2}$. Obviously, K is bounded. Since K is a convex subset of a finite-dimensional space K is totally bounded. Therefore there exist $z^1, \dots, z^k \in K$ with $K \subset \bigcup_{i=1}^k B_{\varepsilon/2}(z^i)$. This implies

$$K = \text{co}(M) \subset \bigcup_{i=1}^k B_{\varepsilon}(z^i).$$

This shows $\text{co}(M)$ is relatively compact. This implies $\overline{\text{co}}(M)$ is compact too. ■

5.9 Some inverse problems

Recovery of signals

A signal is a function in $L_2(\mathbb{R})$. In communications we often deal with signals which are presumed to have no very high frequency components. We make this class of signals precise by use of Fourier transforms and say that a signal $L_2(\mathbb{R})$ is **band limited** with **band** $[-\Omega, \Omega]$ if its Fourier transform \hat{f} vanishes for $|\omega| > \Omega$. With the family of characteristic function χ_a , defined by

$$\chi_a(\omega) := \begin{cases} 1, & |\omega| \leq a \\ 0, & \text{otherwise} \end{cases},$$

a band-limited signal f has the property

$$\hat{f}(\omega) = \chi_{\Omega}(\omega) \hat{f}(\omega), \quad \omega \in \mathbb{R}.$$

We define

$$B_{\Omega} := \{f \in L_2(\mathbb{R}) : \chi_{\Omega} \hat{f} = \hat{f}\},$$

B_{Ω} is a closed subspace of $L_2(\mathbb{R})$.

Let $f \in L_2(\mathbb{R})$ be considered as a signal from which we know the segment $\mathbf{g} := f|_{[-\tau, \tau]}$ ($\tau > 0$). Without any further knowledge about f we cannot reconstruct f from \mathbf{g} . But if we have the information that the signal f is contained in B_Ω for some bounded interval ω then it is possible to reconstruct f from \mathbf{g} . This follows from the fact that f must be an analytic function due to the representations

$$f(t) = \frac{1}{\sqrt{2\pi}} \int_{[-\Omega, \Omega]} \hat{f}(\omega) e^{i\omega t} d\omega, \quad t \in \mathbb{R}. \quad (5.43)$$

If we define

$$D_\tau := \{f \in L_2(\mathbb{R}) : \chi_\tau f = \mathbf{g}\}, \quad (5.44)$$

then our assumptions can be formulated as follows:

$$f \in B_\Omega \cap D_\tau.$$

D_τ is a closed affine subspaces of $L_2(\mathbb{R})$. Therefore, the metric projections of $L_2(\mathbb{R})$ onto D_τ, B_Ω may be considered and used. They are of the following form:

$$\mathcal{D}_\tau : L_2(\mathbb{R}) \ni f \mapsto \chi_\tau f \in D_\tau, \quad (5.45)$$

$$\mathcal{B}_\Omega : L_2(\mathbb{R}) \ni f \mapsto \chi_\Omega \hat{f} \in B_\Omega. \quad (5.46)$$

Set

$$C := \{f \in L_2(\mathbb{R}) : (\chi_\Omega - 1)\hat{f} = \theta\}, \quad D := \{f \in L_2(\mathbb{R}) : \mathcal{D}_\tau f = \mathbf{g}\}.$$

Using these notations the reconstruction may be rewritten as a feasibility problem:

Find a signal $f \in C \cap D$.

Remark 5.16. *A related theme is the problem of phase retrieval of images from data*

$$|\langle \mathbf{a}_j | \mathbf{x} \rangle| = \mathbf{b}_j, \quad j = 1, \dots, m.$$

Consider (see [17])

$$C := \{(\mathbf{x}, \mathbf{z}) \in \mathbb{C}^n \times \mathbb{C}^m : \langle \mathbf{a}_j | \mathbf{x} \rangle = z_j, \quad j = 1, \dots, m\},$$

$$D := \{(\mathbf{x}, \mathbf{z}) \in \mathbb{C}^n \times \mathbb{C}^m : |z_j| = \mathbf{b}_j, \quad j = 1, \dots, m\}.$$

Notice that $\langle \cdot | \cdot \rangle$ is an inner product in \mathbb{C}^n . □

Reconstruction of radiographs

The reconstruction problem of x-ray tomography (x-ray slice imaging) is used to reconstruct the internal structure of a solid object from external measurements. Tomography (($\tau\acute{o}\mu\omicron\sigma$ = slice, $\gamma\rho\alpha\phi\tilde{\iota}\nu$ =aufzeichnen) means to reconstruct a three-dimensional object B from its two-dimensional slices. Objects of interest in x-ray imaging are described by a real-valued function defined on \mathbb{R}^3 , called the attenuation coefficient. The attenuation coefficient quantifies the tendency of an object to absorb or scatter x-rays of a given energy. This function varies from point-to-point within the object and is usually taken to vanishes outside Ω . The attenuation coefficient, like density, is nonnegative. It is useful for medical imaging because different anatomical structures have different attenuation

coefficients. Bone has a much higher attenuation coefficient than soft tissue and different soft tissues have slightly different coefficients. For medical applications it is crucial that normal and cancerous tissues also have slightly different attenuation coefficients.

Using the assumption that x-rays travel along straight lines $I_{\theta,s}$ the x-ray flux is attenuated along this line. Here θ (angle) and s (distance) are parameters identifying the line in a coordinate system. Then a radiograph is a line integral like this:

$$\int_{I_{\theta,s}} f(l) dl$$

where f is the density of the object in the slice. The collection of all integrals of f along the lines in the plane defines a function Rf on the unitsphere in \mathbb{R}^2 times \mathbb{R} . It is called the **Radon transform**. To solve the (inverse) problem of tomography is the task to invert the Radon transform R . J. Radon gave 1917 an inversion formula for R .

To invert the Radon transform in practice there are two main groups of methods in the focus: analytic methods and iterative methods. Analytic methods use mainly properties of the transform R to obtain the density f back from Rf by theoretical considerations. A famous method is the so called **filtered back projection**. Iterative methods solve the linear system of equations which one get by discretization of a finite number of line integrals. As we know from the consideration concerning feasibility such a system may be solved by successive projection on hyperplanes. A famous method is the so called **algebraic reconstruction technique (ART)**. We shall consider this method in the next chapter. As rule, the numerical realization of analytic methods is faster than iterative methods but it needs a lot of physical knowledge. The iterative methods allow more physical and geometrical variations.

A Cauchy problem for an elliptic equation

Let $\Omega \subset \mathbb{R}^2$ be an open bounded and simply connected domain. The problem we consider is the evaluation of the trace of the solution of an elliptic differential equation at that part of the boundary where no data was described, actually at $\partial\Omega \setminus \Gamma$. As a compensation, we observe the Neumann data of the solution on Γ . Such a problem is called a **Cauchy problem**. It is well known that an elliptic Cauchy problem is ill posed mainly due to the lack of continuous dependence of the solution. This can be shown by considering the following Cauchy problem:

$$\begin{aligned} \Delta u(x_1, x_2) &= 0 & , (x_1, x_2) \in \Omega := (0, 1) \times (0, 1) \\ u(x_1, 0) &= 0 & , x_1 \in (0, 1) \\ \frac{\partial}{\partial x_2} u(x_1, 0) &= g_k(x_1) & , x_1 \in (0, 1) \end{aligned} \tag{5.47}$$

where $g_k(x_1) = (\pi k)^{-1} \sin(\pi k x_1)$, $x_1 \in (0, 1)$. It is easily verified that the uniquely determined solution $u = u_k$ is given by

$$u(x_1, x_2) = (\pi k)^{-2} \sin(\pi k x_1) \sinh(\pi k x_2), \quad (x_1, x_2) \in (0, 1) \times (0, 1). \tag{5.48}$$

As we see, the sequence $(g_k)_{k \in \mathbb{N}}$ converges uniformly to zero. Taking the limit $k \rightarrow \infty$ we obtain a Cauchy problem with homogeneous data which admits the trivial solution only. But for every $x_2 > 0$ the solutions u_k oscillate stronger and stronger and became

unbounded as $k \rightarrow \infty$. Consequently, the sequence $(\mathbf{u}_k)_{k \in \mathbb{N}}$ does not converge to zero (in any reasonable topology).

Let us consider the problem (5.48) in a more general setting:

$$\begin{aligned} \mathcal{D}\mathbf{u} &= \boldsymbol{\theta} \quad , \text{ in } \Omega \\ \mathbf{u} &= \mathbf{f} \quad , \text{ in } \Gamma \\ \mathbf{u}_\nu &= \mathbf{g} \quad , \text{ in } \Gamma \end{aligned} \tag{5.49}$$

where \mathcal{D} is a second order elliptic differential operator and \mathbf{f}, \mathbf{g} are given functions on Γ . Let Γ_1 be such that $\partial\Omega = \overline{\Gamma \cup \Gamma_1}$ and $\Gamma \cap \Gamma_1 = \emptyset$. Then we consider the following iteration: Given ϕ_0 on Γ

$$\begin{array}{ll} \mathcal{D}\mathbf{u} = \boldsymbol{\theta} & \text{in } \Omega & \mathcal{D}\mathbf{w} = \boldsymbol{\theta} & \text{in } \Omega \\ \mathbf{u} = \mathbf{f} & \text{in } \Gamma & \mathbf{w} = \boldsymbol{\psi}_k & \text{in } \Gamma_1 \\ \mathbf{u}_\nu = \phi_k & \text{in } \Gamma_1 & \mathbf{w}_\nu = \mathbf{g} & \text{in } \Gamma \end{array}$$

$$\phi_{k+1} := \mathbf{w}_\nu|_{\Gamma_1}.$$

In the iteration step above, two differential equations are to be solved and two trace operators are applied. Actually, we generate two sequences: the first one of Dirichlet traces and the second one of Neumann traces, both defined on Γ_1 . Using trace operators γ_d, γ_n and solution operators for the Dirichlet- and Neumann-problems we conclude that the iteration can be reformulated as

$$\phi_{k+1} = \mathbb{T}(\phi_k) = \mathbb{T}_1(\phi_k) + \mathbf{z}, \quad k \in \mathbb{N}_0. \tag{5.50}$$

with a linear mapping \mathbb{T}_1 and a function \mathbf{z} which depends on the data \mathbf{f}, \mathbf{g} only. From the uniqueness of a Cauchy problem we obtain that a fixed point of \mathbb{T} is a solution of the Cauchy problem. The main difficulty is to find the appropriate spaces of functions defined on the part Γ_1 of the boundary. This can be done (see [40], [8] and [36]) and one obtains that the \mathbb{T}_1 is a linear bounded selfadjoint nonexpansive mapping. An application of the results which we present in Chapter 7 leads to the result that the iteration (5.50) converges to a fixed point of \mathbb{T} .

Online parameter identification

The identification of parameters in partial differential equations arises in a variety of applications. Classical applications are heat conduction problems, population models, seismic explorations and reservoir simulation, computer- and impedance-tomography. The mathematical models which describe these applications contain (function-) parameters which cannot, in general, be measured directly. They must be determined by measuring consequences they produce. Such problems are called inverse problems where the main difficulty to solve identification problems consists in the fact that there is a lack of stability in such problems (ill-posedness).

Here we are concerned with adaptive techniques: determination of the parameter parallel to the process of gathering the data. They have been developed by engineers for control problems governed by ordinary differential equations; see [37], [41], [44]. These techniques can also be used in the case of partial differential equations as it was shown in [3], [11] and [46].

The system we want to consider has the following formulation (in an evolutionary style):

$$D_t z = A(p)z + A_0 z + f(t), t > 0, z(0) = z_0 \quad (5.51)$$

p is the parameter we want to determine from the observation of the state $t \mapsto z(t)$. Later on we weaken this assumption. The Cauchy problem (5.51) is an evolutionary linear system since $A(p), A_0$ are supposed linear operators. For the moment, we assume that the state is defined for all $t > 0$, later on we will consider the case that the process is known on a finite interval $(0, T)$ only. We call the system a **bilinear system** due to the fact that we assume that the mapping $(p, z) \mapsto A(p)z$ is bilinear.

Now, we formulate the equation for the model reference state x as follows:

$$D_t x = Cx + A(q)z(t) + A_0 z(t) - Cz(t) + f(t), t > 0, x(0) = x_0. \quad (5.52)$$

Here C is a mapping with strong ellipticity properties.

In order to find an adaptation rule we exploit the bilinearity of the mapping $(q, v) \mapsto \langle A(q)z(t), v \rangle$:

$$\langle A(h)z(t), v \rangle = \langle B(z(t), v) | h \rangle_{\mathbb{H}}, h \in \mathbb{H}, v \in V, t > 0, \quad (5.53)$$

$$\|B(z(t), \cdot)\| \leq \alpha, t > 0. \quad (5.54)$$

Using this „adjoint“ B we are able to formulate the following adaptation rule:

$$D_t q = -B(z(t), x - z(t)), t > 0, q(0) = q_0. \quad (5.55)$$

Notice that the system (5.52), (5.55) is a nonautonomous system for the evolution of x, q .

We set

$$e := x - z, r := q - p$$

Then we obtain from (5.52), (5.55) the following **error system**:

$$D_t e = Ce + A(r)z(t), t > 0, e(0) = e_0, \quad (5.56)$$

$$D_t r = -B(z(t), e), t > 0, r(0) = r_0. \quad (5.57)$$

In a system notation this reads as follows:

$$D_t \begin{pmatrix} e \\ r \end{pmatrix} = \mathcal{M}(t) \begin{pmatrix} e \\ r \end{pmatrix} \text{ where } \mathcal{M}(t) = \begin{pmatrix} C & A(\cdot)z(t) \\ -B(z(t), \cdot) & \Theta \end{pmatrix}$$

We define

$$L(e, r) := \frac{1}{2} (\|e\|_{\mathbb{H}}^2 + \|r\|_{\mathbb{H}}^2), e \in \mathbb{H}, r \in \mathbb{H}.$$

The main idea for the analysis is to show that L may play the role of a **Ljapunov-function**. A first step is to prove **Output identifiability**:

$$\lim_{t \rightarrow \infty} \|e(t)\|_{\mathbb{H}} = 0. \quad (5.58)$$

To establish parameter convergence, an additional hypothesis is required. Such a hypothesis is an extension of the finite-dimensional notion of *persistence excitation*; see for

instance [43, 47]. Such a hypothesis may be also considered as a *richness*-condition. Under such assumptions **parameter-identifiability** may be proved:

$$\lim_{t \rightarrow \infty} \|r(t)\|_{\mathbb{H}} = 0. \quad (5.59)$$

The model reference adaptive method can be applied also in the case when the state z can be observed in a time interval $[0, T]$ only. This case is similar to the solution of linear systems by using the idea of Kaczmarz which is also known as the ART-procedure in tomography: one uses the state $z(t)$, $t \in [0, T]$, in a cyclic way; see [10, 27, 26, 32]. In this manner we obtain a sequence $(e^n, r^n)_{n \in \mathbb{N}}$ with

$$e^n, r^n : [0, T] \ni t \longmapsto (e^n(t), r^n(t)) \in V \times \mathbb{H}.$$

By analogous reasoning one obtains the results

$$\lim_n \|e^n\|_{\mathbb{H}} = 0 \text{ and } \lim_n \|r^n\|_{\mathbb{H}} = 0;$$

see [46].

The problem statement above has an important drawback: the state has to be completely observed. In the engineering literature this problem is tackled on by introducing a dynamical observer. But this works only well in the „finite-dimensional“ case: the state equation has to be a system of ordinary differential equations. An observer theory for partial differential equations is difficult to realize. Since a partial observation gives occasion to use the projection on the observable part of a state the state z in the model reference system. A promising approach for online identification in the case when a partial observation of the state is available only, is presented by Boiger and Kaltenbacher [13]. Here noisy data are also considered.

5.10 Matrix completion – a non-convex feasibility problem

A partial real matrix is a matrix $A \in \mathbb{R}^{n,n}$ for which entries in certain locations are not known. The problem of matrix completion is the following:

Given a partial matrix find a completion belonging to a specified family of matrices.

This problem can be formulated as a feasibility problem:

Given a partial matrix A find a completion B in a set $C := \bigcap_{i=1}^m C_i$ where C is the intersection of C_1, \dots, C_m . C_1, \dots, C_m are chosen such that the intersection is equal to the set of completions of A with the specified matrix family.

The simplest such case is when C_1 is the set of all completions of A and the intersection C_2, \dots, C_m equals to the desired matrix family. Here we consider the completion problem for the family of Euclidean distance matrices. This is a subclass of the set \mathcal{S}^n of symmetric

matrices in $\mathbb{R}^{n,n}$. Euclidean distance matrices are a useful description of point sets. A typical task is to retrieve the original point configuration. But Euclidean distance matrices may be noisy and/or incomplete.

In this section $\|\cdot\|$ is the Euclidean norm in a space \mathbb{R}^q .

Definition 5.17. Let $D \in \mathcal{S}^n$, $D = (d_{ij})_{i,j=1,\dots,n}$.

(a) D is called a **predistance matrix** if the following hold:

$$d_{ii} = 0, i = 1, \dots, n, a_{ij} \geq 0, i, j = 1, \dots, n, i \neq j.$$

(b) A predistance matrix D is called a **Euclidean distance matrix** if there exist points $x^1, \dots, x^n \in \mathbb{R}^q$ such that

$$\|x^i - x^j\|^2 = d_{ij}, i, j = 1, \dots, n,$$

where here (and in the following) $\|\cdot\|$ is the euclidean norm in \mathbb{R}^n . The points x^1, \dots, x^n are called **generating points** and we say that D is embeddable in \mathbb{R}^q . The lowest number q such that generating points exist \mathbb{R}^q is called the **embedding dimension** $emb(D)$.

□

Notice that generating points of an Euclidean distance matrix are not uniquely determined. This follows from the fact that Euclidean distances are invariant under a translation and a rotation. Moreover, one can easily show that the embedding dimension of a Euclidean distance matrix D satisfies $emb(D) \leq n - 1$.

Suppose we have an Euclidean distance matrix D generated by points $x^1, \dots, x^n \in \mathbb{R}^q$. Then the Gram matrix $G := (\langle x^i | x^j \rangle)_{i,j=1,\dots,n}$ has the representation

$$G = \frac{1}{2}(\mathbf{1}(d^1)^t - D + d^1 \mathbf{1}^t) \text{ where } d^1 \text{ is the first column of } D \quad (5.60)$$

and $\mathbf{1}$ is the vector with entries 1. This leads to the fact that the point set $\{x^1, \dots, x^n\}$ can be regained from G by an eigenvalue decomposition of G .

Euclidean distance matrices have found applications in psychometrics, crystallography, machine learning, wireless sensor networks, acoustics, and more. For a tutorial see the article [20] and [38, 39, 50]. In [39] Euclidean distance matrices are used to study the fractal dimension of attractors of dynamical systems by determining the embedding dimension of associated Euclidean distance matrices. One of the most important application is the study of protein conformation determination. Here methods using the metric projection are successfully applied; see [14, 21]. Now, we follow [14].

To reformulate the completion problem as a feasibility problem we look at the following characterization of a Euclidean distance matrix; see [28].

Theorem 5.18. Let $D \in \mathcal{S}^n$ be a predistance matrix. Then D is a Euclidean distance matrix if and only if the matrix $\hat{D} \in \mathbb{R}^{n-1, n-1}$ in the presentation

$$Q(-D)Q = \begin{pmatrix} \hat{D} & \mathbf{u} \\ \mathbf{u}^t & \alpha \end{pmatrix} \quad (5.61)$$

is positive semi-definite. In this case, the embedding dimension emb of D is $k = rank(\hat{D}) \leq n - 1$.

The matrix Q in (5.61) is the Householder matrix given by

$$Q = I - 2\mathbf{v}\mathbf{v}^t / \|\mathbf{v}\|^2 \text{ where } \mathbf{v} = (1, 1, \dots, 1, 1 + \sqrt{n}) \in \mathbb{R}^n.$$

For a partial predistance matrix $D = (d_{ij})$ consider the associated adjacent matrix $W = (w_{ij}) \in \mathbb{R}^{n,n}$ is defined as follows:

$$w_{ij} = \begin{cases} 1 & \text{if } d_{ij} \text{ is known} \\ 0 & \text{if } d_{ij} \text{ is unknown} \end{cases}$$

We assume that D is embeddable in \mathbb{R}^q . Then we define

$$\begin{aligned} C_1 &:= \{A = (a_{ij}) \in \mathcal{S}^n : a_{ij} \geq 0, a_{ij} = d_{ij} \text{ if } w_{ij} = 1\} \\ C_2 &:= \{A \in \mathcal{S}^n : Q(-A)Q = \begin{pmatrix} \hat{A} & \mathbf{u} \\ \mathbf{u}^t & \alpha \end{pmatrix}, \hat{A} \in \mathcal{S}_+^{n-1}, \mathbf{u} \in \mathbb{R}^{n-1}, \alpha \in \mathbb{R}, \text{rank}(\hat{A}) \leq q\} \end{aligned}$$

Then the completion problem can be reformulated as a feasibility problem

$$\text{Find } A \in C_1 \cap C_2 \tag{5.62}$$

Suppose that $\mathbb{R}^{n,n}$ is endowed with the Frobenius-norm. Clearly C_1 is closed and convex and the projection $P_{C_1}(A)$ of $A \in \mathbb{R}^{n,n}$ can easily be computed. But C_2 is non-convex and does not have the Chebyshev property. Nevertheless, we may consider the projection $P_{C_2}(A)$ of $A \in \mathbb{R}^{n,n}$; we refer to [14] for this computation. Several numerical simulations show that completion of a Euclidean distance matrix is possible by the alternate projection method.

Remark 5.19. *The completion of a Euclidean distance matrix may be done by other approaches: exact completion by graph methods, methods of least squares, approximate completion by semidefinite optimization using interior point methods; see [38] for a short survey.* \square

5.11 Appendix: Variational inequalities

5.12 Conclusions and comments

5.13 Exercises

- 1.) Let $f \in C[a, b]$. What is the best approximation in \mathcal{P}_0 ?
- 2.) Let $f \in C[-1, 1]$, $f(t) := t^3$. Show that $p_0(t) := (1 - t^2)$ is the best approximation of f in the subspace $\text{span}(\{1, t^2\})$ of $C[-1, 1]$.
- 3.) Consider the subspace $\text{span}(\{1, t^2\})$ of $C[-1, 1]$. Show that there exists $f \in C[-1, 1]$ such the best approximation is not uniquely determined.
- 4.) Let $f \in C[a, b]$ and let $p \in \mathcal{P}_n$ with $f(t_i) - p(t_i) = (-1)^i \varepsilon_i$, $\text{sign}(\varepsilon_i) = \text{const}$ at $n + 2$ consecutive points t_0, \dots, t_{n+1} . Show $\text{dist}(f, \mathcal{P}_n) = \min_{i=1, \dots, n+1} |\varepsilon_i|$.

- 5.) Show that the subspace $\text{span}(\{1, t \cos(t), t \sin(t)\})$ of $C[0, \pi]$ satisfies the Haar condition.
- 6.) Let \mathbb{R}^n endowed with l_2 -norm and let

$$\mathbf{U} := \{x = (x_1, \dots, x_n) : \sum_{i=1}^n x_i = 0\}.$$

- (a) Show \mathbf{U} is a linear subspace with $\dim \mathbf{U} = n - 1$.
- (b) Compute $P_{\mathbf{U}}(e^i)$, $i = 1, 2, \dots, n$. (Here: $e^i = (\delta_{ij})$.)
- (c) Compute $P_{\mathbf{U}}(x)$ for all $x \in \mathbb{R}^n$.
- 7.) Let \mathbb{R}^n endowed with the l_2 -norm and let

$$\mathbf{C} := \{x = (x_1, \dots, x_n) : x_1 \leq x_2 \leq \dots \leq x_n\}.$$

- (a) Show that \mathbf{C} is a closed convex cone.
- (b) Show that \mathbf{C} is a Chebyshev set.
- (c) Compute $P_{\mathbf{C}}(e^i)$ for $i = 1, \dots, n$ (Here: $e^i = (\delta_{ij})$.)
- 8.) Show that if a Banach space has a Schauder basis then it is separable.
- 9.) Let $\mathcal{H} = \mathbb{R}^3$ with the Euclidean norm and let

$$\mathbf{H} := \{(x, y, z) \in \mathcal{H} : x + 2y = 0\}.$$

Define on \mathbf{H} the linear functional λ by $\langle \lambda, (x, y, z) \rangle := x - y$. Find $\|\lambda\|$ and extend λ to \mathcal{H} without changing the norm.

- 10.) Let $\mathcal{X} = \mathbb{R}^3$ with the l_1 norm and let

$$\mathbf{U} := \{(x, y, 0) \in \mathcal{X} : x + 2y = 0\}.$$

Define on \mathbf{U} the linear functional λ by $\langle \lambda, (x, y, z) \rangle := x$. Find $\|\lambda\|$ and find two norm-preserving extensions of λ to \mathcal{X} .

- 11.) Let \mathcal{H} be a Hilbert space and let $B : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ be bilinear. Suppose that $B(\cdot, y)$ is continuous for all $y \in \mathcal{H}$ and $B(x, \cdot)$ is continuous for all $x \in \mathcal{H}$. Show that B is continuous, i.e.

$$|B(x, y)| \leq c \|x\| \|y\| \text{ for all } x, y \in \mathcal{H} \text{ for some } c \geq 0.$$

Hint: Uniform boundedness principle.

- 12.) Let \mathcal{X} be a Banach space and let \mathbf{U}, \mathbf{V} be linear subspaces with $\mathcal{X} = \mathbf{U} \oplus \mathbf{V}$. Let $x \in \mathcal{X}, x = u + v, u \in \mathbf{U}, v \in \mathbf{V}$; define $Px := v$. Then we have a mapping $P : \mathcal{X} \rightarrow \mathbf{V}$. Then

- (a) P is linear and $P \circ P = P$.
- (b) \mathbf{U}, \mathbf{V} are closed if and only if P is continuous.

Bibliography

- [1] S. Al-Homidan. Hybrid methods for solving the educational testing problem. *J. of Computational and Applied Mathematics*, 91:31–45, 1998.
- [2] F. Albiac, L.L. Ansorena, S.J. Dillworth, and D. Kutzarova. Existence and uniqueness of greedy bases in banach spaces. <https://arxiv.org/abs/1505.08119>, x:1–28, 2015.
- [3] H.W. Alt, K.-H. Hoffmann, and J. Sprekels. A numerical procedure of solve certain identification problems. *Intern. Ser. Numer. Math.*, 68:11–43, 1984.
- [4] J. Baumeister. Über die Extremaleigenschaft nichtlinearer Splines. *Numer. Math.*, 25:433–455, 1976.
- [5] J. Baumeister. Variationsprobleme in Orliczräumen und Splines. *Manuscripta Math.*, 20:29–49, 1977.
- [6] J. Baumeister. *Stable Solution of Inverse Problems*. Vieweg, Braunschweig, 1987.
- [7] J. Baumeister. Konvexe Analysis, 2014. Skriptum WiSe 2014/15, Goethe-Universität Frankfurt/Main.
- [8] J. Baumeister and A. Leitão. On iterative methods for solving ill-posed problems modeled by PDE's. *Journal of Inverse and Ill-Posed Problems*, 9:13–29, 2001.
- [9] J. Baumeister and L. Schumaker. Nonlinear classes of splines and variational problems. *J. of Approximation Theory*, 18:63–73, 1976.
- [10] J. Baumeister, B. Kaltenbacher and A. Leitão. On Levenberg-Marquardt-Kaczmarz iterative methods for solving systems of nonlinear equations. *Inverse Problems and Imaging*, 4:335–350, 2010.
- [11] J. Baumeister, W. Scondo, M.A. Demetriou and I.G. Rosen. On-line parameter estimation for infinite dimensional dynamical systems. *SIAM Journal on Control and Optimization*, 35:678 – 713, 1997.
- [12] D. Bengs. Greedy Approximation im Raum der Funktionen von beschränkter Variation und Anwendungen auf die Bildverarbeitung. Master's thesis, Goethe-Universität, 2012. Fachbereich Informatik und Mathematik.
- [13] R. Boiger and B. Kaltenbacher. A online parameter identification method for time dependent partial differential equations. *Inverse Problems*, page 28 pages, 2015.
- [14] J.M. Borwein and M.K. Tam. Reflection methods for inverse problems with applications to protein conformation determination. pages 1–18.

- [15] E.W. Cheney and A. Goldstein. Proximity maps for convex sets. *Proc. Amer. Math. Soc.*, 10:448–450, 1959.
- [16] M.T. Chu and J.W. Wright. The educational testing problem revisited. *IMA J. Numer. Anal.*, 15:141–160, 1995.
- [17] P.L. Combettes. The convex feasibility problem in image recovery. *Advances in Imaging and Electron Physics*, 95:155–270, 1996.
- [18] S.J. Dillworth, N.J. Kalton, and D. Kutzarova. On the existence of almost greedy bases in banach spaces. *XXX*, pages xx–xx, 2012.
- [19] S.J. Dillworth, D. Freman, E. Odell and T. Schlumprecht. Greedy bases for Besov spaces. *Constr. Approximation*, pages xx–xx, 2010.
- [20] I. Dokmanic, R. Parhizkar, J. Ranieri, and M. Vetterli. Euclidean distance matrices. arXiv:1502.07541v2, 2015.
- [21] Q. Dong and Z. Wu. A geometric build-up algorithm for solving the molecular distance geometric problem with sparse distance data. *J. Global Optim.*, 26:321–333, 2003.
- [22] P. Enflo. A counterexample to the approximation problem in Banach spaces. *Acta Mathematica*, 130:309–317, 1973.
- [23] R. Fletcher. A nonlinear programming problem in statistics (educational testing). *SIAM J. Sci. Stat. Comput.*, 2:257–267, 1981.
- [24] R. Fletcher. Semi-definite matrix constraints. *SIAM J. Control and Optimization*, 23:493–513, 1985.
- [25] W. Glunt. An alternating projections method for certain linear problems in a hilbert space. *IMA Journal of Numerical Analysis*, 15:291–305, 1995.
- [26] M. Haltmeier, R. Kowar, A. Leitao, and O. Scherzer. Kaczmarz methods for regularizing nonlinear ill-posed equations. II. Applications. *Inverse Probl. Imaging*, 1(3):507–523, 2007.
- [27] M. Haltmeier, A. Leitao, and O. Scherzer. Kaczmarz methods for regularizing nonlinear ill-posed equations. I. convergence analysis. *Inverse Probl. Imaging*, 1(2):289–298, 2007.
- [28] T.L. Hayden and J. Wells. Approximation by matrices positive semidefinite on a subspace. *Linear Algebra and its Applications*, 109:115–130, 1988.
- [29] N. Higham. Computing a nearest symmetric positive semi-definite matrix. *Linear Algebra and Appl.*, 103:103–118, 1988.
- [30] N. Higham. Computing the nearest correlation matrix – a problem from finance, 2002.
- [31] R. Holmes. *A Course on Optimization and Best Approximation*. Springer, 1971.
- [32] S. Kaczmarz. Approximate solution of systems of linear equations. *Internat. J. Control*, 57:1269–1271, 1993.
- [33] A. Kirsch. *An Introduction to the Mathematical Theory of Inverse Problems*. Springer, New York, 1996.

- [34] G. Koethe. *Topological linear spaces*. Springer, Berlin, 1960.
- [35] S.V. Konyagin and V.N. Temlyakov. A remark on greedy approximation in banach spaces. *East J. Approx.*, 5:365–379, 1999.
- [36] Kozlov, V.G. Maz’ya, and A.V. Fomin. An iterative method for solving the cauchy problem for elliptic equations. *Comput. Maths. Phys.*, 31:45–52, 1991.
- [37] G. Kreisselmeier. Adaptive observers with exponential rate of convergence. *IEEE Trans. Autom. Contr.*, 22:509 – 535, 1977.
- [38] Miriam Kreth. Matrixergänzungsprobleme. Master’s thesis, Goethe-Universität, 2004. Fachbereich Informatik und Mathematik.
- [39] Miriam Kreth. *Distanzmatrizen, erzeugende Punkte und Einbettungsdimension*. Dissertation, Goethe Universität, 2008.
- [40] A. Leitao. An iterative method for solving elliptic Cauchy problems. *Numer. Func. Anal. Optim.*, 21:715–742, 2000.
- [41] P.M. Lion. Rapid identification of linear and nonlinear systems. *AIAA J.*, 5:1835 – 1841, 1987.
- [42] B.J. McCartin. Theory of exponential splines. *J. of Approximation theory*, 66:1–23, 1991.
- [43] K.S. Narendra and A.M. Annaswamy. Persistent excitation in adaptive systems. *Int. J. Control*, 45:127 – 160, 1987.
- [44] K.S. Narendra and P. Kudva. Stable adaptive schemes for system identification and control i,ii. *IEEE SMC-4*, 4:542 – 560, 1974.
- [45] F. Natterer. *The Mathematics of Computerized Tomography*. John Wiley, New York, 1986.
- [46] W. Scondo. Ein Modellabgleichungsverfahren zur adaptiven Parameteridentifikation in Evolutionsgleichungen. Master’s thesis, Fachbereich Mathematik, Goethe-Universität, Frankfurt/Main, 1987.
- [47] N. Shimkin and A. Feuer. Persistency of excitation in continuous-time systems. *Systems & Control Letters*, 9:225–233, 1987.
- [48] V.N. Temlyakov. The best m-term approximation and greedy algorithms. *Advances in Computational Mathematics*, 8:249–265, 1998.
- [49] V.N. Temlyakov. Greedy approximation. *Acta Mathematica*, 17:1069–1083, 2008.
- [50] Annika Wacker. Das Matrix-Ergänzungsproblem und die Berechnung einer besten Approximation einer Korrelationsmatrix. Master’s thesis, Goethe Universität, 2016. Fachbereich Informatik und Mathematik.
- [51] D. Werner. *Funktionalanalysis*. Springer, 2002.