## **Optimization Methods**

with Applications in Statistics

Lutz Dümbgen University of Bern

July 11, 2025

# **Bibliography**

- [1] A. BEN-TAL and A. NEMIROVSKI (2001). Lectures on Modern Convex Optimization. Analysis, Algorithms, and Engineering Applications. SIAM, Philadelphia
- [2] R.J. BERAN and L. DÜMBGEN (2010). Least Squares and Shrinkage Estimation under Bimonotonicity Constraints. *Statistics and Computing* 20(2), 177–189.
- [3] CORMEN, T.H., C.E. LEISERSON and R.L. RIVEST (1990). Introduction to Algorithms. M.I.T. Press
- [4] K. DEIMLING (1985). Nonlinear Functional Analysis. Springer, Berlin Heidelberg
- [5] R. FLETCHER (1987). Practical Methods of Optimization (2nd edition). Wiley, New York
- [6] P.E. GILL, W. MURRAY and M.H. WRIGHT (1981). Practical Optimization. Academic Press
- [7] U. GRENANDER (1956), On the Theory of Mortality Measurement, Part II. Skandinavisk Aktuarietidskrift **39**, 125–153.
- [8] F. JARRE und J. STOER (2004). Optimierung. Springer, Berlin Heidelberg
- [9] K. LANGE (1999). Numerical Analysis for Statisticians. Springer, New York
- [10] G. OPFER (1994). Numerische Mathematik für Anfänger. Vieweg, Braunschweig Wiesbaden.
- [11] T. ROBERTSON and P. WALTMAN (1968). On Estimating Monotone Parameters. Annals of Mathematical Statistics 39, 1030–1039
- [12] T. ROBERTSON, F.T. WRIGHT and R.L. DYKSTRA (1988). Order Restricted Statistical Inference. Wiley, New York

Acknowledgements. Sebastian Arnold, David Ginsbourger, Micha Grütter, Dirk Klingbiel, Jürg Krähenbühl, Christof Mahnig, Alexandre Mösching, Kaspar Rufibach, Adrian Schmucker, Blendi Shala and Christof Strähl gave numerous hints to different versions of these lecture notes – thanks a lot to all of them!

# Contents

1	1 Univariate Procedures				
	1.1	Bisection Methods		11	
	1.2	2 Newton's Method (Univariate Case)			
	1.3	Golden Section Search			
2	Con	vex Sets	5	25	
	2.1	1 Convex Sets and Cones			
	2.2	2.2 Metric Projections and Separating Hyperplanes			
		2.2.1	Metric projections	28	
		2.2.2	Support functions and half spaces	32	
		2.2.3	Separating hyperplanes	33	
		2.2.4	Polar sets and cones	35	
	2.3	.3 Extremal Points			
	2.4 Convex Polyhedra				
		2.4.1	Linear programs	43	
		2.4.2	Doubly stochastic matrices and the Hoffmann–Wielandt inequalities	47	
		2.4.3	Polyhedral cones	50	
3	Con	vex Fur	actions	53	
	3.1	Conve	x Functions on the Real Line	53	
	3.2 Convex Functions on Linear Spaces		x Functions on Linear Spaces	59	
	3.3	3 Directional Derivatives and Minimizers			
	3.4 Smoothness Properties			68	
		3.4.1	Local Lipschitz continuity	68	
		3.4.2	Subdifferentials and smoothness	72	

		3.4.3	Jensen's inequality and Bregman divergences	74				
	3.5	Lower	Semicontinuity, Minimizers and Convex Conjugates	77				
		3.5.1	Lower semicontinuity	77				
		3.5.2	Convexity, lower semicontinuity and affine functions	78				
		3.5.3	Existence of minimizers and coercivity	80				
		3.5.4	Convex conjugates	83				
4	Mul	tivariat	e Optimization	87				
	4.1	Newto	n's Method	87				
	4.2	Minim	vization Problems	94				
		4.2.1	A general criterion for convergence	95				
		4.2.2	Gradient, Newton and quasi-Newton procedures	98				
		4.2.3	Examples for the candidate step function $\Delta$	99				
		4.2.4	Examples for the step size function $\lambda$	100				
		4.2.5	Performance in connection with quasi-Newton methods	106				
5	Con	straineo	d Optimization	109				
	5.1	Lagrar	nge Multipliers	109				
		5.1.1	The general principle	109				
		5.1.2	Examples for Lagrange's Method	110				
		5.1.3	Justification of Lagrange's method	121				
		5.1.4	Lagrange duality	125				
	5.2	Specia	l Algorithms	127				
		5.2.1	Iterative algorithms	128				
		5.2.2	Active set methods	132				
		5.2.3	Isotonic least squares regression	138				
		5.2.4	The pool-adjacent-violators algorithm (PAVA)	140				
6	Conjugate Gradients 15							
	6.1	The Ta	ısk	151				
	6.2	The G	radient Method	152				
	6.3	Conjug	gate Directions	155				
	6.4	The Co	onjugate Gradient (CG) Algorithm	156				
	6.5	Bound	ing the running time and approximation error	159				

	6.6	Minimizing a Smooth Convex Function	168			
7	7 Dynamic Programming					
	7.1	Dykstra's Algorithm	171			
	7.2	Alignment of Sequences	173			
	7.3	Bimonotone Regression	177			

# Introduction

In this course we are dealing with various optimization methods which are used frequently in statistics and other fields. The optimization problems we'll encounter may be divided roughly into three categories:

(a) Minimization problems. A function  $f : \mathcal{X} \to \overline{\mathbb{R}} = [-\infty, \infty]$  has to be minimized. This concerns both the minimal value,

$$\inf_{\mathcal{X}} f = \inf_{x \in \mathcal{X}} f(x),$$

as well as the set of minimizers,

$$\underset{\mathcal{X}}{\operatorname{arg\,min}} f = \underset{x \in \mathcal{X}}{\operatorname{arg\,min}} f(x) := \Big\{ x \in \mathcal{X} : f(x) = \underset{\mathcal{X}}{\operatorname{inf}} f \Big\}.$$

(b) Finding roots. For a mapping  $f : \mathcal{X} \to \mathcal{Y}$  and a point  $b \in \mathcal{Y}$  we want to determine the set

$$f^{-1}(b) = \{x \in \mathcal{X} : f(x) = b\}$$

or at least one of its elements.

(c) Fixed point problems. For a mapping  $f : \mathcal{X} \to \mathcal{Y}$  with  $\mathcal{X} \subset \mathcal{Y}$  we want to determine its set of fixed points,

$$FP(f) := \{ x \in \mathcal{X} : f(x) = x \}.$$

Quite often one translates minimization problems into root finding or fixed point problems.

## Chapter 1

# **Univariate Procedures**

In this chapter we consider functions  $f : \mathcal{X} \to \mathbb{R}$  on a real interval  $\mathcal{X} \subset \mathbb{R}$ .

## **1.1 Bisection Methods**

Suppose that  $f : [a, b] \to \mathbb{R}$  is a continuous function such that  $f(a) \le 0 \le f(b)$ . According to the intermediate value theorem for continuous functions, there exists at least one point  $x_* \in [a, b]$  with  $f(x_*) = 0$ . This fact is used in the algorithm in Table 1.1. It returns a pair  $(x_a, x_b)$  of points in [a, b] such that  $0 \le x_b - x_a \le \delta$  and  $f(x_*) = 0$  for some  $x_* \in [x_a, x_b]$ . It terminates after

$$\left\lceil \log_2\left(\frac{b-a}{\delta}\right) \right\rceil$$

executions of the while-loop, because after  $n \ge 0$  executions,

$$x_b - x_a = (b - a)2^{-n}$$

Sometimes one would also like to ensure that  $|f(x)| \leq \delta$  for all  $x \in [x_a, x_b]$ . To do so, we now assume that f is also isotonic (non-decreasing) and modify the algorithm from Table 1.1 as described in Table 1.2. In case of a differentiable function f one can verify easily that the

```
\begin{array}{l} \textbf{Algorithm } (x_a, x_b) \leftarrow \textbf{Bisection}(f, a, b, \delta) \\ (x_a, x_b) \leftarrow (a, b) \\ \textbf{while } x_b - x_a > \delta \ \textbf{do} \\ x_o \leftarrow (x_a + x_b)/2 \\ \textbf{if } f(x_o) \leq 0 \ \textbf{then} \\ x_a \leftarrow x_o \\ \textbf{else} \\ x_b \leftarrow x_o \\ \textbf{end if} \\ \textbf{end while.} \end{array}
```

Table 1.1: Bisection method I.

Algorithm 
$$(x_a, x_b) \leftarrow Bisection 2(f, a, b, \delta)$$
  
 $(x_a, x_b) \leftarrow (a, b)$   
 $(f_a, f_b) \leftarrow (f(a), f(b))$   
while  $x_b - x_a > \delta$  or  $f_b - f_a > \delta$  do  
 $x_o \leftarrow (x_a + x_b)/2$   
 $f_o \leftarrow f(x_o)$   
if  $f_o \le 0$  then  
 $x_a \leftarrow x_o$   
 $f_a \leftarrow f_o$   
else  
 $x_b \leftarrow x_o$   
 $f_b \leftarrow f_o$   
end if  
end while.

Table 1.2: Bisection method II.

while-loop is executed no more than

$$\left\lceil \log_2 \left( \frac{b-a}{\delta} \max \left\{ 1, \sup_{[a,b]} f' \right\} \right) \right\rceil$$

times.

#### Application to exact confidence bounds

Let Y be a random variable with values in  $\mathbb{Z}$  and distribution function  $F_{\theta_*}$ , where  $(F_{\theta})_{\theta \in \Theta}$  is a given family of distribution functions, and  $\theta_*$  is an unknown parameter in  $\Theta$ . Thus

$$\mathbb{P}(Y \le c) = F_{\theta_*}(c).$$

One can easily show that

$$\mathbb{P}(F_{\theta_*}(Y) > \alpha) \ge 1 - \alpha$$
 and  $\mathbb{P}(F_{\theta_*}(Y - 1) < 1 - \alpha) \ge 1 - \alpha$ 

for any fixed  $\alpha \in (0, 1)$ . Thus both

$$C_{\alpha}^{\text{left}}(Y) := \left\{ \theta \in \Theta : F_{\theta}(Y) > \alpha \right\}$$

and

$$C^{\mathrm{right}}_{\alpha}(Y) \ := \ \left\{ \theta \in \Theta : F_{\theta}(Y-1) < 1-\alpha \right\}$$

are  $(1 - \alpha)$ -confidence regions for  $\theta_*$ . The set  $C_{\alpha}^{\text{left}}(Y)$  consists of all parameters  $\theta$  such that Y is not "suspiciously small" in the sense that the left-sided p-value  $F_{\theta}(Y)$  is larger than  $\alpha$ . Analogously,  $C_{\alpha}^{\text{right}}(Y)$  consists of all parameters  $\theta$  such that Y is not "suspiciously large" in the sense that the right-sided p-value  $1 - F_{\theta}(Y - 1)$  is larger than  $\alpha$ .

Suppose that  $\Theta$  is an interval of real numbers. For any integer  $c \in \Theta \ni \theta \mapsto F_{\theta}(c)$  be continuous and antitonic (non-increasing). In this case  $C_{\alpha}^{\text{left}}(Y)$  yields an upper and  $C_{\alpha}^{\text{right}}(Y)$  yields a lower confidence bound for  $\theta_*$ .

More precisely, let  $F_{\theta}(Y)$  be continuous and strictly antitonic in  $\theta \in \Theta$  with limits

$$\lim_{\theta \to \inf(\Theta)} F_{\theta}(Y) = 1 \quad \text{and} \quad \lim_{\theta \to \sup(\Theta)} F_{\theta}(Y) = 0.$$

Then

$$C_{\alpha}^{\text{left}}(Y) = \left\{ \theta \in \Theta : \theta < b_{\alpha}(Y) \right\}$$

with the unique number  $b_{\alpha}(Y) \in \Theta$  such that

$$F_{b_{\alpha}(Y)}(Y) = \alpha.$$

Thus we obtain an exact upper confidence bound which may be computed via a bisection procedure.

Analogously let  $F_{\theta}(Y-1)$  be continuous and strictly antitonic in  $\theta \in \Theta$  with limits

$$\lim_{\theta \to \inf(\Theta)} F_{\theta}(Y-1) = 1 \text{ and } \lim_{\theta \to \sup(\Theta)} F_{\theta}(Y-1) = 0.$$

Then

$$C_{\alpha}^{\mathrm{right}}(Y) = \left\{ \theta \in \Theta : \theta > a_{\alpha}(Y) \right\}$$

with the unique number  $a_{\alpha}(Y) \in \Theta$  such that

$$F_{a_{\alpha}(Y)}(Y-1) = 1 - \alpha.$$

**Example 1.1** (Binomial parameters). Let Y follow a binomial distribution  $Bin(n, \theta_*)$  with given  $n \in \mathbb{N}$  and unknown parameter  $\theta_* \in \Theta = [0, 1]$ . Then

$$F_{\theta}(c) = \sum_{k=0}^{c} {n \choose k} \theta^{k} (1-\theta)^{n-k}$$

for arbitrary  $c \in \{0, 1, ..., n\}$ . As a function of  $\theta$ , this is a polynomial of order n and thus infinitely often differentiable. In case of  $0 \le c < n$  one can even show that  $F_{\theta}(c)$  is strictly antitonic in  $\theta \in [0, 1]$  with limits  $F_0(c) = 1$  and  $F_1(c) = 0$ . Consequently,

$$C_{\alpha}^{\text{left}}(Y) = \begin{cases} [0,1] & \text{if } Y = n, \\ [0,b_{\alpha}(Y)) & \text{if } Y < n, \text{ where } F_{b_{\alpha}(Y)}(Y) = \alpha, \end{cases}$$

and

$$C_{\alpha}^{\text{right}}(Y) = \begin{cases} [0,1] & \text{if } Y = 0, \\ (a_{\alpha}(Y),1] & \text{if } Y > 0, \text{ where } F_{a_{\alpha}(Y)}(Y-1) = 1 - \alpha. \end{cases}$$

For the computation of both confidence bounds  $a_{\alpha}(Y)$  and  $b_{\alpha}(Y)$  one may employ the algorithm in Table 1.3 which returns for given numbers  $c \in \{0, 1, ..., n-1\}, \gamma \in (0, 1)$  and  $\delta > 0$  two numbers  $x_a, x_b$  with the following properties:

$$0 \le x_a < x_b \le 1, \ x_b - x_a \le \delta,$$

and

$$F_{x_a}(c) \geq \gamma \geq F_{x_b}(c), \ F_{x_a}(c) - F_{x_b}(c) \leq \delta.$$

It requires an auxiliary function  $FBino(\cdot, n, \theta)$  for the computation of the distribution function of  $Bin(n, \theta)$ .

$$\begin{array}{l} \textbf{Algorithm} \ (x_a, x_b) \leftarrow \textbf{BinoCB}(n, c, \gamma, \delta) \\ (x_a, x_b) \leftarrow (0, 1) \\ (F_a, F_b) \leftarrow (1, 0) \\ \textbf{while} \ x_b - x_a > \delta \ \textbf{or} \ F_a - F_b > \delta \ \textbf{do} \\ x_o \leftarrow (x_a + x_b)/2 \\ F_o \leftarrow FBino(c, n, x_o) \\ \textbf{if} \ F_o \ge \gamma \ \textbf{then} \\ x_a \leftarrow x_o \\ F_a \leftarrow F_o \\ \textbf{else} \\ x_b \leftarrow x_o \\ F_b \leftarrow F_o \\ \textbf{end if} \\ \textbf{end while.} \end{array}$$

Table 1.3: Computation of confidence bounds for a binomial parameter.

**Example 1.2** (Poisson parameters). Let Y follow a Poisson distribution  $Poiss(\theta_*)$  with unknown parameter  $\theta_* \in \Theta = [0, \infty)$ . Then

$$F_{\theta}(c) = \exp(-\theta) \sum_{k=0}^{c} \frac{\theta^{k}}{k!}$$

for arbitrary  $c \in \mathbb{N}_0$ . As a function of  $\theta$ , this is continuous and strictly antitonic with limits  $F_0(c) = 1$  and  $F_{\infty}(c) = 0$ . Thus

$$C^{\mathrm{left}}_{\alpha}(Y) \ = \ [0, b_{\alpha}(Y)) \quad \text{with} \ F_{b_{\alpha}(Y)}(Y) = \alpha,$$

and

$$C^{\mathrm{right}}_{\alpha}(Y) \;=\; \begin{cases} [0,\infty) & \text{if } Y=0, \\ (a_{\alpha}(Y),\infty) & \text{if } Y>0, \text{ where } F_{a_{\alpha}(Y)}(Y-1)=1-\alpha. \end{cases}$$

For the computation of both confidence bounds  $a_{\alpha}(Y)$  and  $b_{\alpha}(Y)$  one has to modify algorithm **BinoCB** in two respects. On the one hand we replace the function  $FBino(\cdot, n, \theta)$  with the distribution function  $FPoiss(\cdot, \theta)$  of  $Poiss(\theta)$ . Moreover, at first one has to find a starting interval for the bisection search, see Table 1.4.

**Exercise 1.3.** Let  $w_0, w_1, w_2, \ldots$  be nonnegative weights such that  $w_0 > 0$  and

$$C(\theta) := \sum_{k=0}^{\infty} w_k \theta^k < \infty \text{ for all } \theta \ge 0.$$

With these weights we define a family of distributions on  $\mathbb{N}_0$ : For  $\theta \ge 0$  and  $c \in \mathbb{N}_0$  let

$$F_{\theta}(c) := C(\theta)^{-1} \sum_{k=0}^{c} w_k \theta^k.$$

(a) Show that  $\theta \mapsto F_{\theta}(c)$  is continuous and antitonic on  $[0, \infty)$  with  $F_0(c) = 1$ . Show that it is even strictly antitonic with limit  $F_{\infty}(c) = 0$ , provided that  $w_k > 0$  for at least one index k > c.

```
Algorithm (x_a, x_b) \leftarrow \mathbf{PoissCB}(c, \gamma, \delta)
(x_a, x_b) \leftarrow (0, 1)
(F_a, F_b) \leftarrow (1, \texttt{FPoiss}(c, 1))
while F_b > \gamma do
         (x_a, x_b) \leftarrow (x_b, 2x_b)
          (F_a, F_b) \leftarrow (F_b, \texttt{FPoiss}(c, x_b))
end while
while x_b - x_a > \delta or F_a - F_b > \delta do
          x_o \leftarrow (x_a + x_b)/2
          F_o \leftarrow \texttt{FPoiss}(c, x_o)
         if F_o \geq \gamma then
                   x_a \leftarrow x_o
                    F_a \leftarrow F_o
         else
                   x_b \leftarrow x_o
                    F_b \leftarrow F_o
          end if
end while.
```

Table 1.4: Computation of confidence bounds for a Poisson parameter.

(b) This result applies immediately to Poisson distributions, where  $w_k = 1/k!$ . How could you apply it to binomial distributions? How would you modify the result to be applicable to negative binomial distributions?

### **1.2** Newton's Method (Univariate Case)

Now we consider a differentiable function  $f : \mathcal{X} \to \mathbb{R}$  on an open interval  $\mathcal{X} \subset \mathbb{R}$  with strictly positive derivative f'. Moreover we assume that  $\inf_{\mathcal{X}} f < 0 < \sup_{\mathcal{X}} f$ , hence there exists a unique  $x_* \in f^{-1}(0)$ .

To find or approximate this point  $x_*$  we consider an arbitrary first candidate  $x_o \in \mathcal{X}$ . Now we approximate f by an affine function:

$$x \mapsto f(x_o) + f'(x_o)(x - x_o).$$

This auxiliary function attains the value 0 at

$$\psi(x_o) := x_o - \frac{f(x_o)}{f'(x_o)}.$$

Note that  $\psi(x_o) = x_o$  if and only if  $f(x_o) = 0$ , whence  $x_o = x_*$ . Otherwise we hope that  $\psi(x_o)$  is closer to  $x_*$  than  $x_o$ .

Now we would like to iterate this mapping  $\psi$ . That means, for a starting value  $x_0 \in \mathcal{X}$  one computes inductively

$$x_n := \psi(x_{n-1})$$
 for  $n = 1, 2, 3, \dots$ 

When doing so we hope that

- (i) this sequence is well-defined in the sense that all points  $x_n$  are contained in  $\mathcal{X}$ ,
- (ii) it converges to  $x_*$ .

**Local convergence.** The next theorem implies that the sequence  $(x_n)_{n=0}^{\infty}$  does have these properties, provided that the starting value  $x_0$  is sufficiently close to  $x_*$ .

Theorem 1.4 (Local convergence of Newton's method).

(a) Let f' be continuous at  $x_*$ . Then

$$\lim_{x \to x_*} \frac{\psi(x) - x_*}{x - x_*} = 0.$$

(If  $x = x_*$ , then  $\psi(x) = x_*$ , and we interpret 0/0 as 0.)

(b) Let f be twice differentiable, and let f'' be continuous at  $x_*$ . Then

$$\lim_{x \to x_*} \frac{\psi(x) - x_*}{(x - x_*)^2} = \frac{f''(x_*)}{2f'(x_*)}$$

For the sequence  $(x_n)_{n=0}^{\infty}$  defined earlier, Theorem 1.4 (a) has the following consequences: For any  $\epsilon \in (0, 1)$  there exists a  $\delta(\epsilon) > 0$  such that  $[x_* \pm \delta(\epsilon)] \subset \mathcal{X}$  and

$$|\psi(x) - x_*| \leq \epsilon |x - x_*|$$
 for  $x \in [x_* \pm \delta(\epsilon)]$ .

In particular, if  $|x_0 - x_*| \leq \delta(\epsilon)$ , then  $|x_n - x_*| \leq |x_0 - x_*|\epsilon^n$  for all  $n \geq 0$ . More generally, if  $|x_{n(\epsilon)} - x_*| \leq \delta(\epsilon)$  for some  $n(\epsilon) \in \mathbb{N}_0$ , then  $|x_n - x_*| \leq C(\epsilon)\epsilon^n$  for all  $n \geq 0$ , where  $C(\epsilon) := \max_{m \leq n(\epsilon)} \epsilon^{-m} |x_m - x_*|$ .

This consideration shows that

(1.1) 
$$\lim_{n \to \infty} |x_n - x_*| = 0 \quad \text{and} \quad \lim_{n \to \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|} = 0,$$

provided that  $x_0$  is sufficiently close to  $x_*$ . We'll show later how to avoid the latter restriction. Property (1.1) of  $(x_n)_{n\geq 0}$  is called "super-linear convergence (to  $x_*$ )".

For starting values  $x_0$  sufficiently close to  $x_*$ , Theorem 1.4 (b) yields even the more precise statement that

(1.2) 
$$\lim_{n \to \infty} |x_n - x_*| = 0 \quad \text{und} \quad \lim_{n \to \infty} \frac{|x_{n+1} - x_*|}{|x_n - x_*|^2} = \frac{|f''(x_*)|}{2f'(x_*)}.$$

This property of  $(x_n)_{n=0}^{\infty}$  is called "quadratic convergence (to  $x_*$ )".

**Proof of Theorem 1.4.** For real numbers a, b let  $a \wedge b$  and  $a \vee b$  denote their minimum and maximum, respectively. According to the mean value theorem, for  $x \in \mathcal{X} \setminus \{x_*\}$ , there exists a point  $\xi(x) \in (x \wedge x_*, x \vee x_*)$  such that

$$f(x) = f(x) - f(x_*) = f'(\xi(x))(x - x_*),$$

whence

$$\frac{\psi(x) - x_*}{x - x_*} = 1 - \frac{f(x)}{f'(x)(x - x_*)} = 1 - \frac{f'(\xi(x))(x - x_*)}{f'(x)(x - x_*)} = 1 - \frac{f'(\xi(x))}{f'(x)}.$$

$$\lim_{x \to x_*} \frac{\psi(x) - x_*}{x - x_*} = 1 - \lim_{x \to x_*} \frac{f'(\xi(x))}{f'(x)} = 1 - \frac{f'(x_*)}{f'(x_*)} = 0.$$

This proves part (a).

For part (b) we use Taylor's formula. For  $x \in \mathcal{X} \setminus \{x_*\}$  and  $h := x - x_*$ ,

$$f(x) = f(x_*) + f'(x_*)h + f''(\xi(x))h^2/2 = f'(x_*)h + f''(\xi(x))h^2/2$$

and

$$f'(x) = f'(x_*) + f''(\eta(x))h$$

with suitable points  $\xi(x)$ ,  $\eta(x)$  in  $(x \wedge x_*, x \vee x_*)$ . Thus,

$$\frac{\psi(x) - x_*}{(x - x_*)^2} = \frac{h - f(x)/f'(x)}{h^2}$$

$$= \frac{f'(x)h - f(x)}{f'(x)h^2}$$

$$= \frac{f'(x_*)h + f''(\eta(x))h^2 - f'(x_*)h - f''(\xi(x))h^2/2}{(f'(x_*) + f''(\eta(x))h)h^2}$$

$$= \frac{f''(\eta(x)) - f''(\xi(x))/2}{f'(x_*) + f''(\eta(x))h}$$

$$\to \frac{f''(x_*) - f''(x_*)/2}{f'(x_*)} = \frac{f''(x_*)}{2f'(x_*)}$$

as  $x \to x_*$ , i.e.  $h \to 0$ .

**Example 1.5**. Let us illustrate the potential problems with starting values  $x_0$  which are too far from  $x_*$  with the function  $f : \mathbb{R} \to \mathbb{R}$ ,  $f(x) := x(1+x^2)^{-1/2}$ . This function is strictly increasing with  $f'(x) = (1+x^2)^{-3/2}$ , and the limits are  $f(\pm \infty) = \pm 1$ . Here  $x_* = 0$ , and

$$\psi(x) = x - x(1 + x^2) = -x^3$$

Hence the Newton sequence  $(x_n)_{n\geq 0}$  converges if and only if  $|x_0| < 1$ . Precisely, one can show by induction that

$$x_n = (-1)^n x_0^{3^n}$$

for all  $n \ge 0$ .

**Global convergence.** Under the additional assumption that f is convex and  $x_0 \ge x_*$ , Newton's method yields always a sequence with limit  $x_*$ . This is a consequence of the following result:

**Theorem 1.6** (Global convergence of Newton's method). Suppose that f is convex, that means, f' is isotonic. Then  $\psi(x) \ge x_*$  whenever  $x \le x_*$ . For  $x > x_*$ ,

$$0 \leq \frac{\psi(x) - x_*}{x - x_*} \leq 1 - \frac{f'(x_* +)}{f'(x)}$$

where  $f'(x_* +)$  denotes the right limit of f' at  $x_*$ .

Under the assumptions of Theorem 1.6, the sequence  $(x_n)_{n\geq 0}$  is monotone decreasing with limit  $x_*$ , provided that  $x_0 \geq x_*$ . The monotonicity is even strikt, unless  $x_n = x_*$  for some  $n \geq 0$ . But  $\psi(x) = x_*$  for some  $x > x_*$  would imply that  $f' \equiv f'(x_* +)$  on  $(x_*, x)$ .

The restriction  $x_0 \ge x_*$  is unnecessary if  $\mathcal{X} = (a, \infty)$  for some  $a \in [-\infty, \infty)$ . For then  $x_1 = \psi(x_0) \ge x_*$ , and  $(x_n)_{n=1}^{\infty}$  is monotone decreasing with limit  $x_*$ .

**Example 1.7.** A classical application of Newton's procedure is the computation of square roots. For any number  $\gamma > 0$  we consider the function

$$(0,\infty) \ni x \mapsto f(x) := x^2 - \gamma.$$

This function is infinitely often differentiable, strictly isotonic and convex, and  $f^{-1}(0) = \{\sqrt{\gamma}\}$ . Here

$$\psi(x) = x - \frac{x^2 - \gamma}{2x} = \frac{x + \gamma/x}{2}$$

and iterating this mapping yields always a sequence converging quadratically to  $\sqrt{\gamma}$ . Indeed,

$$\frac{\psi(x) - \sqrt{\gamma}}{x - \sqrt{\gamma}} = \frac{x - \sqrt{\gamma}}{2x} \begin{cases} < 1/2 \\ > 0 & \text{if } x > \sqrt{\gamma} \end{cases}$$

and

$$\frac{\psi(x) - \sqrt{\gamma}}{(x - \sqrt{\gamma})^2} = \frac{1}{2x}$$

for any x > 0.

**Proof of Theorem 1.6.** For  $x < x_*$  there exists a point  $\xi(x) \in (x, x_*)$  such that

$$\psi(x) = x - \frac{f'(\xi(x))(x - x_*)}{f'(x)} = x_* + (x - x_*) \left(1 - \frac{f'(\xi(x))}{f'(x)}\right) \ge x_*,$$

because  $0 < f'(x) \le f'(\xi(x))$ .

For  $x > x_*$  there exists a point  $\xi(x) \in (x_*, x)$  with

$$\frac{\psi(x) - x_*}{x - x_*} = 1 - \frac{f'(\xi(x))}{f'(x)} \in \left[0, 1 - \frac{f'(x_* +)}{f'(x)}\right],$$

because  $f'(x_* +) \le f'(\xi(x)) \le f'(x)$ .

**Exercise 1.8.** Consider the function  $f : (0, \infty) \to \mathbb{R}$ ,  $f(x) := x^k - \gamma$ , with given constants k > 1 and  $\gamma > 0$ .

(a) Show that f is strictly isotonic and convex with unique zero at  $x_* = \gamma^{1/k}$ .

(b) Show that the corresponding algorithmic mapping  $\psi(x) = x - f(x)/f'(x)$  for Newton's method has the following properties:

$$h(x) := \frac{\psi(x) - x_*}{x - x_*}$$

is isotonic on  $(x_*,\infty)$  with limits  $h(x_*) = 0$  and  $h(\infty) = 1 - 1/k$ .

**Exercise 1.9.** For  $k \in \mathbb{N}$  let  $F : [0, \infty) \to \mathbb{R}$  be given by

$$F(x) := 1 - \exp(-x) \sum_{j=0}^{k} x^j / j!.$$

(a) Show that F is continuously differentiable and strictly isotonic on  $[0, \infty)$  with F(0) = 0 and limit  $F(\infty) = 1$ .

(b) Describe Newton's procedure for the (approximate) solution of the equation  $F(x) = \gamma$ , where  $\gamma$  is any given number in (0, 1).

(c) Specify a starting value  $x_0 > 0$  such that Newton's method yields for arbitrary  $\gamma \in (0, 1)$  a convergent sequence.

**Stopping criteria.** So far we analyzed the algorithmic mapping  $\psi$  underlying Newton's method. For explicit implementations we also need a stopping criterion involving some upper bound  $\delta$  for the approximation error. Since we are looking for a zero of f, we could run the algorithm until  $|f(x_n)| \leq \delta$ . Sometimes one combines this criterium with the requirement that  $|x_n - x_{n-1}| \leq \delta$ . Suppose we want to guarantee that  $|x_n - x_*| \leq \delta$ . If we can show that  $f'(x) \geq \kappa > 0$  for all x in an interval containing  $x_*$  and all points  $x_n$ , then it follows from

$$\frac{f(x_n)}{x_n - x_*} = \frac{f(x_n) - f(x_*)}{x_n - x_*} = f'(\xi_n)$$

with suitable  $\xi_n \in (x_n \wedge x_*, x_n \vee x_*)$  that

$$|x_n - x_*| \leq \frac{|f(x_n)|}{\kappa}.$$

Thus the stopping criterion  $|f(x_n)| \leq \delta \kappa$  guarantees that  $|x_n - x_*| \leq \delta$ .

A simpler rule of thumb which is often implemented relies on the fact that in case of (1.1) also

$$\lim_{n \to \infty} \frac{x_n - x_*}{x_n - x_{n-1}} = 0,$$

see the next exercise. Consequently, if we iterate the mapping  $\psi$  until  $|x_n - x_{n-1}|$  is smaller than a (very) small number  $\delta$ , chances are high that  $|x_n - x_*| \leq \delta$ , too.

**Exercise 1.10** (Rates of convergence). Let  $(\mathcal{X}, d)$  be a metric space and  $(x_n)_{n=0}^{\infty}$  a sequence in  $\mathcal{X}$  with limit  $x_* \in \mathcal{X}$ , that is,  $\lim_{n\to\infty} d(x_n, x_*) = 0$ .

(a) Suppose that

$$\lim_{n \to \infty} \frac{d(x_{n+1}, x_*)}{d(x_n, x_*)} = 0$$

(with the conventions that 0/0 := 0 and  $a/0 := \infty$  for a > 0). Show that

$$\lim_{n \to \infty} \frac{d(x_n, x_*)}{d(x_n, x_{n-1})} = 0$$

Show also that for any  $\epsilon \in (0,1)$  there exists a  $K = K(\epsilon) \in (0,\infty)$  such that

$$d(x_n, x_*) \leq K\epsilon^n$$
 for all  $n \geq 0$ 

(b) Suppose that

$$\frac{d(x_{n+1}, x_*)}{d(x_n, x_*)^2} \le C < \infty \quad \text{for all } n \ge 0$$

Show that for each  $\epsilon \in (0, 1)$  and  $n, k \in \mathbb{N}_0$ ,

$$d(x_{n+k}, x_*) \leq \epsilon^{(2^k)}/C \quad \text{if } d(x_n, x_*) \leq \epsilon/C.$$

**Exercise 1.11.** For a given constant c > 1, let  $f(x) := -\log(1-x) - cx$ .

- (a) Show that there is a unique point  $x_* \in (0, 1)$  such that  $f(x_*) = 0$ .
- (b) Describe explicit programs to approximate  $x_*$  up to a given error  $\delta \in (0,1)$  by means of a
- (b.1) bisection method,
- (b.2) Newton method.

Your descriptions should specify how to determine a starting interval or starting point, respectively, and provide a stopping criterion such that the approximation  $x_n$  of  $x_*$  satisfies  $|x_n - x_*| \le \delta$ .

**Exercise 1.12**. For a given parameter c > 1, let  $f(x) := e^x - 1 - cx$ .

- (a) Show that there is a unique point  $x_* > 0$  such that  $f(x_*) = 0$ .
- (b) Show that  $\log c < x_* < 2 \log c$ .

(c) Show that  $f'(x_*) = h(x_*) := e^{x_*} - (e^{x_*} - 1)/x_*$ . Then show that h is strictly positive and strictly isotonic on  $(0, \infty)$ , and deduce from that the inequality  $f'(x_*) > c - (c-1)/\log c > 0$ .

(d) Describe an explicit Newton procedure to approximate  $x_*$ . For a given threshold  $\delta > 0$ , your procedure should return a number  $x_{**} > 0$  such that  $|x_{**} - x_*| \le \delta$  and  $|f(x_{**})| \le \delta$ .

**Exercise 1.13**. Suppose that f is convex on  $\mathcal{X}$  but f'(x) is concave in  $x \ge x_*$ . Show that

$$0 \le rac{\psi(x) - x_*}{x - \psi(x)} < 1$$
 for any  $x > x_*$ .

For a Newton sequence  $(x_n)_{n\geq 0}$  with starting value  $x_0 \geq x_*$  this implies that

$$0 \leq x_n - x_* \leq x_{n-1} - x_n$$
 for arbitrary  $n \geq 1$ .

**Exercise 1.14.** Consider independent random variables  $X_1, X_2, \ldots, X_n \in [0, 1]$  with density function

$$f_{\theta}(x) := c(\theta)^{-1} \exp(\theta x), \quad x \in [0, 1],$$

where  $\theta \in \mathbb{R}$  is some unknown parameter and  $c(\theta) > 0$  some norming constant.

(a) Compute  $c(\theta)$  and verify that

$$\mu(\theta) := \mathbb{E}_{\theta}(X_1) = \frac{d}{d\theta} \log c(\theta).$$

(b) Show that  $\mu$  is differentiable with

$$\mu'(\theta) = \operatorname{Var}_{\theta}(X_1) > 0$$

and verify that its limits are  $\mu(-\infty) = 0$  and  $\mu(\infty) = 1$ .

- (c) Show that  $\mu(-\theta) = 1 \mu(\theta)$  for all  $\theta \in \mathbb{R}$ .
- (d) Show that  $\mu(\theta)$  is convex in  $\theta \in (-\infty, 0]$  and concave in  $\theta \in [0, \infty)$ .
- (e) Show that the log-likelihood function

$$\theta \mapsto \sum_{i=1}^n \log f_\theta(X_i)$$

has a unique maximizer  $\hat{\theta} \in \mathbb{R}$  which is given by

$$\bar{X} = \mu(\theta)$$

with the sample mean  $\bar{X} = n^{-1} \sum_{i=1}^{n} X_i$ .

(f) Describe and implement an explicit procedure to compute  $\hat{\theta}$ .

*Hint:* For small values of  $|\theta|$  (e.g.  $|\theta| < 0.01$ ) the function  $\mu(\theta)$  may be numerically instable. (Why?) Use a suitable Taylor approximation of  $\mu$  to avoid such problems.

### **1.3 Golden Section Search**

We end this chapter with a simple procedure to minimize a "bathtub-shaped" function  $f : [a, b] \rightarrow (-\infty, \infty]$ . Precisely, suppose that there exists a unique minimizer  $x_*$  of f. Further let f be strictly antitonic on  $[a, x_*]$  and strictly isotonic on  $[x_*, b]$ . Our goal is to find (an approximation of)  $x_*$  without referring to any derivative of f.

To localize the minimizer  $x_*$  we rely on the following simple consideration: Let  $a \le r < s \le b$ . Then it follows from  $f(r) \le f(s)$  that  $x_* < s$ . For if  $x_* \ge s$ , then f(r) > f(s), because f is strictly antitonic on  $[a, x_*] \supset [r, s]$ . Analogously it follows from  $f(r) \ge f(s)$  that  $x_* > r$ .

This consideration gives rise to an iterative procedure: Let  $a \le q < r < s < t \le b$ , where it is known that  $x_* \in [q, t]$ . To fulfill the latter requirement, we start with q = a and t = b. In case of  $f(r) \le f(s)$  we know that  $x_* \in [q, s]$ , so we may replace the quadruple (q, r, s, t) with  $(q, \dot{r}, r, s)$ , where  $q < \dot{r} < r$ . In case of f(r) > f(s) it is clear that  $x_* \in [r, t]$ , so we replace (q, r, s, t) with  $(r, s, \dot{s}, t)$ , where  $s < \dot{s} < t$ , see figure 1.1.

For the explicit choice of the new points  $\dot{r}$  and  $\dot{s}$  there are various possibilities. For instance one could take  $\dot{r} = (q+r)/2$  and  $\dot{s} = (s+t)/2$ . Here is another proposal: We want to ensure that the three quadrupels (q, r, s, t),  $(q, \dot{r}, r, s)$  and  $(r, s, \dot{s}, t)$  are identical up to affine transformations. To this end we consider the ratios

$$B := \frac{r-q}{t-q}$$
 and  $C := \frac{s-q}{t-q}$ ,

so 0 < B < C < 1 and

$$r = (1 - B)q + Bt, \quad s = (1 - C)q + Ct$$

The quadrupels  $(q, \dot{r}, r, s)$  and (q, r, s, t) have the same 'shape' if and only if

$$\frac{\dot{r}-q}{r-q} = \frac{r-q}{s-q} = \frac{B}{C}$$
 and  $\frac{r-q}{s-q} = \frac{s-q}{t-q} = C$ ,



Figure 1.1: Partitions for Golden Section Search.

so

$$B = C^2$$
 and  $\dot{r} = q + C(r - q) = (1 - C)q + Cr.$ 

The quadrupels  $(r, s, \dot{s}, t)$  and (q, r, s, t) have the same shape if and only if

$$\frac{t-\dot{s}}{t-s} = \frac{t-s}{t-r} = \frac{1-C}{1-B}$$
 and  $\frac{t-s}{t-r} = \frac{t-r}{t-q} = 1-B$ ,

so

$$1 - C = (1 - B)^2$$
 and  $\dot{s} = t - (1 - B)(t - s) = (1 - B)s + Bt$ 

Plugging in  $B = C^2$  into  $1 - C = (1 - B)^2$  yields the equation

$$C^4 - 2C^2 + C = 0.$$

Obviously, C = 0 and C = 1 are two irrelevant solutions of this equation. Thus we write

$$C^{4} - 2C^{2} + C = C(C^{3} - 2C + 1) = C(C - 1)(C^{2} + C - 1)$$

and determine a solution  $C \in (0, 1)$  of the equation

$$0 = C^{2} + C - 1 = (C + 1/2)^{2} - 5/4.$$

This leads to

$$C = \frac{\sqrt{5}-1}{2} = \frac{2}{\sqrt{5}+1} \approx 0.618.$$

Then the corresponding value for  $B = C^2 = 1 - C$  equals

$$B = \frac{3 - \sqrt{5}}{2} = \frac{2}{3 + \sqrt{5}} \approx 0.382.$$

In figure 1.1 we used these particular ratios B and C. Since B + C = 1, we may rewrite

$$r = Cq + Bt$$
,  $s = Bq + Ct$ ,  $\dot{r} = Bq + Cr$ ,  $\dot{s} = Cs + Bt$ .

These considerations lead to the algorithm in Table 1.5. It yields an interval [q, t] of length at most  $C\delta$  containing the minimizer  $x_*$ . Note that it terminates after at most

$$\left\lceil \log_{1/C} \left( \frac{b-a}{\delta} \right) \right\rceil$$

repetitions of the while-loop.

**Algorithm**  $(q, t) \leftarrow$  **GoldenSection** $(f, a, b, \delta)$  $(B, C) \leftarrow (0.382, 0.618)$  $(q, r, s, t) \leftarrow (a, Ca + Bb, Ba + Cb, b)$  $(f_r, f_s) \leftarrow (f(r), f(s))$ while  $t - q > \delta$  do if  $f_r \leq f_s$  then  $(r, s, t) \leftarrow (Bq + Cr, r, s)$  $(f_r, f_s) \leftarrow (f(r), f_r)$ else  $(q, r, s) \leftarrow (r, s, Cs + Bt)$  $(f_r, f_s) \leftarrow (f_s, f(s))$ end if end while if  $f_r \leq f_s$  then  $t \leftarrow s$ else  $q \leftarrow r$ end if.

Table 1.5: Golden section search.

Exercise 1.15. Implement the golden section search to approximate the minimizer of the function

$$x \mapsto f(x) := (1-x)\log(1-x) + x - x^2$$

on [0, 1] with given precision  $\delta > 0$ . (Here  $0 \log(0) := 0$ .) Why is golden section search justified here? What is your result in case of  $\delta = 0.0001$ ?

**Exercise 1.16**. What can be said about the algorithm for golden section search if we assume only that f is

- (a) continuous,
- (b) differentiable,
- (c) continuous and convex

on [*a*, *b*]?

**Exercise 1.17.** Suppose one replaces algorithm GoldenSection $(f, a, b, \delta)$  with the algorithm in Table 1.6.

(a) Which ratios

$$(B,C) := \left(\frac{r-q}{t-q}, \frac{s-q}{t-q}\right)$$

can occur while running this algorithm?

(b) In what sense is the new algorithm worse than golden section search?

```
Algorithm (q, t) \leftarrow SilverSection(f, a, b, \delta)
(q, r, s, t) \leftarrow (a, 2a/3 + b/3, a/3 + 2b/3, b)
(f_r, f_s) \leftarrow (f(r), f(s))
while t - q > \delta do
        if f_r \leq f_s then
                (r,s,t) \leftarrow ((q+r)/2,r,s)
                 (f_r, f_s) \leftarrow (f(r), f_r)
        else
                 (q,r,s) \leftarrow (r,s,(s+t)/2)
                 (f_r, f_s) \leftarrow (f_s, f(s))
        end if
end while
if f_r \leq f_s then
        t \gets s
else
        q \leftarrow r
end if.
```

Table 1.6: 'Silver section' search.

## Chapter 2

# **Convex Sets**

In the current and next chapter we cover some basic concepts and results of convex analysis. In what follows we define and analyze properties of subsets of a real vector space V.

## 2.1 Convex Sets and Cones

**Definition 2.1** (Convex sets). A subset C of V is called *convex*, if for arbitrary points  $x, y \in C$  their connecting line segment is contained in C as well:

$$\left\{ (1-t)\boldsymbol{x} + t\boldsymbol{y} : t \in [0,1] \right\} \subset \boldsymbol{C}.$$

Via induction this implies that

$$\sum_{i=1}^n \lambda_i oldsymbol{x}_i \ \in \ oldsymbol{C}$$

for arbitrary  $n \in \mathbb{N}$ ,  $x_1, \ldots, x_n \in C$  and  $\lambda_1, \ldots, \lambda_n \ge 0$  with  $\sum_{i=1}^n \lambda_i = 1$ . Such a point  $\sum_{i=1}^n \lambda_i x_i$  is called a "convex combination of the points  $x_1, \ldots, x_n$  (with weights  $\lambda_1, \ldots, \lambda_n$ )".

**Example 2.2.** Suppose that  $(V, \|\cdot\|)$  is a normed vector space. Then any closed ball  $B(x_o, \delta) := \{x \in V : \|x - x_o\| \le \delta\}$  with center  $x_o \in V$  and radius  $\delta \ge 0$  is a convex set. Indeed, if  $x, y \in B(x_o, \delta)$  and  $\lambda \in [0, 1]$ , then

$$\begin{aligned} \left\| (1-\lambda)\boldsymbol{x} + \lambda\boldsymbol{y} - \boldsymbol{x}_o \right\| &= \left\| (1-\lambda)(\boldsymbol{x} - \boldsymbol{x}_o) + \lambda(\boldsymbol{y} - \boldsymbol{x}_o) \right\| \\ &\leq \left\| (1-\lambda)(\boldsymbol{x} - \boldsymbol{x}_o) \right\| + \left\| \lambda(\boldsymbol{y} - \boldsymbol{x}_o) \right\| \\ &= (1-\lambda) \|\boldsymbol{x} - \boldsymbol{x}_o\| + \lambda \|\boldsymbol{y} - \boldsymbol{x}_o\| \\ &\leq (1-\lambda)\delta + \lambda\delta = \delta. \end{aligned}$$

With the same arguments one can show that any open ball is a convex set too.

**Definition 2.3** (Convex cone). A subset C of V is called a *convex cone*, if for arbitrary points  $x, y \in C$  the set of their linear combinations with nonnegative coefficients is contained in C as well:

$$ig\{ \lambda oldsymbol{x} + \mu oldsymbol{y} : \lambda, \mu \geq 0 ig\} \ \subset \ oldsymbol{C}$$

In other words, C is convex, and for arbitrary  $x \in C$  the set of its nonnegative multiples is contained in C:

$$\{\lambda \boldsymbol{x}: \lambda \geq 0\} \subset \boldsymbol{C}.$$

Again one can show via induction that

$$\sum_{i=1}^n \lambda_i oldsymbol{x}_i \ \in \ oldsymbol{C}$$

for arbitrary  $n \in \mathbb{N}$ ,  $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n \in \boldsymbol{C}$  and  $\lambda_1, \ldots, \lambda_n \geq 0$ .

**Convex hulls etc.** Now let M be an arbitrary nonvoid subset of V. Then one can show easily that

$$\operatorname{conv}(\boldsymbol{M}) := \left\{ \sum_{i=1}^n \lambda_i \boldsymbol{x}_i : n \in \mathbb{N}; \, \boldsymbol{x}_1, \dots, \boldsymbol{x}_n \in \boldsymbol{M}; \, \lambda_1, \dots, \lambda_n \ge 0; \, \sum_{i=1}^n \lambda_i = 1 \right\}$$

is the smallest convex set containing M. We call this set the convex hull of M. Moreover,

$$\operatorname{cone}(\boldsymbol{M}) := \left\{ \sum_{i=1}^n \lambda_i \boldsymbol{x}_i : n \in \mathbb{N}; \, \boldsymbol{x}_1, \dots, \boldsymbol{x}_n \in \boldsymbol{M}; \, \lambda_1, \dots, \lambda_n \ge 0 \right\}$$

is the smallest convex cone containing M. We call it the *convex cone spanned by* M. In case of a finite-dimensional space V, these representations may be simplified as follows:

Lemma 2.4 (Carathéodory<sup>1</sup>, 1907). Let  $d := \dim(V)$  be finite. Then for arbitrary subsets M of V,

$$\operatorname{conv}(\boldsymbol{M}) = \left\{ \sum_{i=1}^{d+1} \lambda_i \boldsymbol{x}_i : \boldsymbol{x}_1, \dots, \boldsymbol{x}_{d+1} \in \boldsymbol{M}; \, \lambda_1, \dots, \lambda_{d+1} \ge 0; \, \sum_{i=1}^{d+1} \lambda_i = 1 \right\},\\ \operatorname{cone}(\boldsymbol{M}) = \left\{ \sum_{i=1}^{d} \lambda_i \boldsymbol{x}_i : \boldsymbol{x}_1, \dots, \boldsymbol{x}_d \in \boldsymbol{M}; \, \lambda_1, \dots, \lambda_d \ge 0 \right\}.$$

**Proof of Lemma 2.4.** We start with the representation of cone(M). Let  $\boldsymbol{x} = \sum_{i=1}^{n} \lambda_i \boldsymbol{x}_i$  be a point in cone(M), where  $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n \in M, \lambda_1, \ldots, \lambda_n > 0$  and n > d. Then the vectors  $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$  are linearly dependent. That means,

$$\sum_{i=1}^n \mu_i x_i = \mathbf{0}$$

for suitable numbers  $\mu_1, \ldots, \mu_n$  such that  $\mu_i \neq 0$  for at least one index *i*. Without loss of generatily let  $\mu_i < 0$  for at least one index *i*. Then for arbitrary  $t \ge 0$ ,

$$\boldsymbol{x} = \sum_{i=1}^{n} (\lambda_i + t \mu_i) \boldsymbol{x}_i.$$

<sup>&</sup>lt;sup>1</sup>Constantin Carathéodory (1873–1950): Greek mathematician who spent most of his academic life in Germany with significant contributions to measure theory and other fields.

If we choose

$$t := \min_{i: \mu_i < 0} \frac{\lambda_i}{|\mu_i|},$$

then all coefficients  $\lambda_i := \lambda_i + t\mu_i$  are nonnegative, and at least one of them equals zero. Thus we may represent x as a linear combination of n - 1 vectors in cone(M) with nonnegative coefficients. We may iterate this reduction until x is a linear combination of d vectors in cone(M) with nonnegative coefficients.

As to conv(M), note that the equations  $x = \sum_{i=1}^{n} \lambda_i x_i$  and  $\sum_{i=1}^{n} \lambda_i = 1$  with  $x_i \in M$  and  $\lambda_i > 0$  may be combined to one equation

$$(\boldsymbol{x},1) = \sum_{i=1}^n \lambda_i(\boldsymbol{x}_i,1)$$

involving points in the (d+1)-dimensional linear space  $V \times \mathbb{R}$ . If n > d+1, we may imitate the arguments for cone(M) to show that (x, 1) is a linear combination of n-1 points in  $M \times \{1\}$  with non-negative weights. Again iteration of this argument shows that (x, 1) is a linear combination of d+1 points in  $M \times \{1\}$  with non-negative weights. In other words, x is a convex combination of d+1 points in M.

**Exercise 2.5**. Determine the convex hull of  $M \subset \mathbb{R} \times \mathbb{R}$  in the following cases:

- (a)  $M = \{(x, 1/|x|) : x \in \mathbb{R}, x \neq 0\}.$
- **(b)**  $M = \{(x, 1/x) : x \in \mathbb{R}, x \neq 0\}.$

**Exercise 2.6** (Linear mappings and convexity). Let  $L : V \to W$  be a linear mapping from V into another real vector space W.

- (a) Show that L(C) is convex for any convex set  $C \subset V$ .
- (b) Show that  $L^{-1}(D)$  is convex for any convex set  $D \subset W$ .

**Exercise 2.7.** Let  $(V, \|\cdot\|)$  be a real normed vector space and M be a nonvoid subset of V. Prove or falsify the following claims:

- (a) If M is open, then conv(M) is open, too.
- (b) If M is closed, then conv(M) is closed, too.
- (c) If M is compact and  $\dim(V) < \infty$ , then  $\operatorname{conv}(M)$  is compact, too.

**Exercise 2.8**. Let *C* be a convex subset of a real normed vector space  $(V, \|\cdot\|)$ . Show that its closure  $\overline{C}$  is convex, too.

**Exercise 2.9** (Minkowski operations). For  $n \in \mathbb{N}$ , nonvoid sets  $C_1, C_2, \ldots, C_n \subset V$  and real numbers  $\lambda_1, \lambda_2, \ldots, \lambda_n$  let

$$\sum_{i=1}^n \lambda_i \boldsymbol{C}_i := \Big\{ \sum_{i=1}^n \lambda_i \boldsymbol{x}_i : \boldsymbol{x}_i \in \boldsymbol{C}_i \text{ for } 1 \le i \le n \Big\}.$$

Especially for  $x \in V$  and  $C \subset V$  we set  $x + C := \{x\} + C$ .

(a) Show that  $\sum_{i=1}^{n} \lambda_i C_i$  is convex if all sets  $C_i$  are convex.

(b) Show that for a nonvoid convex set  $C \subset V$  and nonnegative numbers  $\lambda_1, \lambda_2, \ldots, \lambda_n$ ,

$$\sum_{i=1}^{n} (\lambda_i C) = \left( \sum_{i=1}^{n} \lambda_i \right) C.$$

**Exercise 2.10**. Let M be a subset of  $\mathbb{R}^d$ .

(a) Suppose that M is finite. Show that cone(M) is closed. (Hint: Consider first the situation that the elements of M are linearly independent. Then argue as in the proof of Lemma 2.4.)

(b) Find an example of a compact set  $M \subset \mathbb{R}^2$  such that  $\operatorname{cone}(M)$  is not closed.

## 2.2 Metric Projections and Separating Hyperplanes

In this section let  $(V, \langle \cdot, \cdot \rangle, \| \cdot \|)$  be a real Hilbert space, for instance,  $V = \mathbb{R}^d$  with the standard inner product  $\langle x, y \rangle := x^\top y$ . The closed unit ball is denoted by B,

$$B = \{ x \in V : ||x|| \le 1 \}.$$

Throughout this section let C be a nonvoid convex and closed subset of V.

#### 2.2.1 Metric projections

The next theorem shows that for any point  $x \in V$  there exists a unique closest point  $\Pi x = \Pi_C x$ in C, that means,

$$\|\boldsymbol{x} - \Pi \boldsymbol{x}\| = d(\boldsymbol{x}, \boldsymbol{C}) := \inf_{\boldsymbol{y} \in \boldsymbol{C}} \|\boldsymbol{x} - \boldsymbol{y}\|.$$

The corresponding mapping  $\Pi : V \to C$  is called the *metric projection of* V onto C. It is Lipschitz continuous with constant one:

Theorem 2.11 (Metric projection onto a closed convex set).

(a) To each point  $x \in V$  there exists a unique point  $\Pi x \in C$  such that

$$\|x - \Pi x\| = \min_{y \in C} \|x - y\|.$$

(b) A point  $y_o \in C$  equals  $\Pi x$  if and only if

$$\langle \boldsymbol{x} - \boldsymbol{y}_o, \boldsymbol{y} - \boldsymbol{y}_o \rangle \leq 0$$
 for all  $\boldsymbol{y} \in \boldsymbol{C}$ .

(c) For arbitrary  $x_1, x_2 \in V$ ,

$$\|\Pi x_1 - \Pi x_2\| \le \|x_1 - x_2\|.$$

In case of closed convex cones or closed linear subspaces, the characterization in part (b) of Theorem 2.11 may be simplified somewhat, and the metric projection  $\Pi$  has additional properties: **Corollary 2.12** (Metric projection onto a closed convex cone). Let C be a closed convex cone. Then a point  $y_o \in C$  equals  $\Pi x$  if and only if

$$\langle \boldsymbol{x} - \boldsymbol{y}_o, \boldsymbol{y}_o \rangle = 0$$
 and  $\langle \boldsymbol{x} - \boldsymbol{y}_o, \boldsymbol{y} \rangle \leq 0$  for all  $\boldsymbol{y} \in \boldsymbol{C}$ .

Moreover,  $\Pi(\lambda x) = \lambda \Pi x$  for arbitrary  $x \in V$  and  $\lambda \ge 0$ .

**Corollary 2.13** (Metric projection onto a closed linear space). Let C be a closed linear subspace of V. Then a point  $y_o \in C$  equals  $\Pi x$  if and only if

$$\langle \boldsymbol{x} - \boldsymbol{y}_o, \boldsymbol{y} \rangle = 0$$
 for all  $\boldsymbol{y} \in \boldsymbol{C}$ .

Moreover,  $\Pi: V \to C$  is a linear mapping, the so-called orthogonal projection from V onto C.

**Proof of Theorem 2.11.** Proof of part (a). We first show that the distance between two points  $y_1, y_2 \in C$  has to be rather small if both  $||x - y_1||$  and  $||x - y_2||$  are close to d(x, C). To this end we use the parallelogram identity,

$$\|a + b\|^2 + \|a - b\|^2 = 2\|a\|^2 + 2\|b\|^2$$
 für  $a, b \in V$ ,

which is a consequence of the equations  $\|a \pm b\|^2 = \|a\|^2 \pm 2\langle a, b \rangle + \|b\|^2$ . This identity is applied to  $a = x - y_1$  and  $b = x - y_2$ , and we utilize the fact that due to convexity of C,

$$\boldsymbol{y}_o := 2^{-1}(\boldsymbol{y}_1 + \boldsymbol{y}_2)$$

belongs to C, too. Thus

$$\begin{split} \|\boldsymbol{y}_{1} - \boldsymbol{y}_{2}\|^{2} &= \|(\boldsymbol{x} - \boldsymbol{y}_{1}) - (\boldsymbol{x} - \boldsymbol{y}_{2})\|^{2} \\ &= 2\|\boldsymbol{x} - \boldsymbol{y}_{1}\|^{2} + 2\|\boldsymbol{x} - \boldsymbol{y}_{2}\|^{2} - \|(\boldsymbol{x} - \boldsymbol{y}_{1}) + (\boldsymbol{x} - \boldsymbol{y}_{2})\|^{2} \\ &= 2\|\boldsymbol{x} - \boldsymbol{y}_{1}\|^{2} + 2\|\boldsymbol{x} - \boldsymbol{y}_{2}\|^{2} - 4\|\boldsymbol{x} - \boldsymbol{y}_{o}\|^{2} \\ &\leq 2\|\boldsymbol{x} - \boldsymbol{y}_{1}\|^{2} + 2\|\boldsymbol{x} - \boldsymbol{y}_{2}\|^{2} - 4d(\boldsymbol{x}, \boldsymbol{C})^{2}, \end{split}$$

whence

$$\|\boldsymbol{y}_1 - \boldsymbol{y}_2\|^2 \le 2(\|\boldsymbol{x} - \boldsymbol{y}_1\|^2 - d(\boldsymbol{x}, \boldsymbol{C})^2) + 2(\|\boldsymbol{x} - \boldsymbol{y}_2\|^2 - d(\boldsymbol{x}, \boldsymbol{C})^2).$$

This inequality shows already that there is at most one point  $y \in C$  such that ||x - y|| = d(x, C). It also implies existence of such a point. To see this, let  $(y_n)_{n=1}^{\infty}$  be a sequence in C such that

$$\lim_{n\to\infty} \|\boldsymbol{x}-\boldsymbol{y}_n\| = d(\boldsymbol{x},\boldsymbol{C}).$$

For  $N \in \mathbb{N}$ ,

$$\begin{split} \sup_{m,n\geq N} \|\boldsymbol{y}_m - \boldsymbol{y}_n\|^2 &\leq 4 \Big( \sup_{n\geq N} \|\boldsymbol{x} - \boldsymbol{y}_n\|^2 - d(\boldsymbol{x}, \boldsymbol{C})^2 \Big) \\ &\rightarrow 0 \quad \text{as } N \rightarrow \infty. \end{split}$$

Hence  $(\boldsymbol{y}_n)_{n=1}^{\infty}$  is a Cauchy sequence in  $\boldsymbol{V}$ . Since  $(\boldsymbol{V}, \|\cdot\|)$  is complete and  $\boldsymbol{C}$  is closed, this sequence has a limit  $\boldsymbol{y} \in \boldsymbol{C}$ , and  $\|\boldsymbol{x} - \boldsymbol{y}\| = d(\boldsymbol{x}, \boldsymbol{C})$ .

Proof of part (b). For any two points  $y_o, y \in C$  and  $t \in [0, 1]$ , the point  $y_t := (1 - t)y_o + ty$ belongs to C, too, and

$$\|\boldsymbol{x} - \boldsymbol{y}_t\|^2 = \|\boldsymbol{x} - \boldsymbol{y}_o - t(\boldsymbol{y} - \boldsymbol{y}_o)\|^2 = \|\boldsymbol{x} - \boldsymbol{y}_o\|^2 - 2t\langle \boldsymbol{x} - \boldsymbol{y}_o, \boldsymbol{y} - \boldsymbol{y}_o\rangle + t^2 \|\boldsymbol{y} - \boldsymbol{y}_o\|^2.$$

The right hand side is greater than or equal to  $\| \boldsymbol{x} - \boldsymbol{y}_o \|^2$  for arbitrary  $t \in [0, 1]$  if and only if

$$\langle oldsymbol{x} - oldsymbol{y}_o, oldsymbol{y} - oldsymbol{y}_o 
angle \ \leq \ 0.$$

Proof of part (c). According to part (b),

$$egin{aligned} \|m{x}_1 - m{x}_2\|^2 &= ig\| (\Pi m{x}_1 - \Pi m{x}_2) + (m{x}_1 - \Pi m{x}_1 - m{x}_2 + \Pi m{x}_2) ig\|^2 \ &\geq \|\Pi m{x}_1 - \Pi m{x}_2\|^2 \ &- 2 \langle m{x}_1 - \Pi m{x}_1, \Pi m{x}_2 - \Pi m{x}_1 
angle - 2 \langle m{x}_2 - \Pi m{x}_2, \Pi m{x}_1 - \Pi m{x}_2 
angle \ &\geq \|\Pi m{x}_1 - \Pi m{x}_2\|^2. \end{aligned}$$

Hence  $\Pi$  is Lipschitz continuous with constant 1.

**Proof of Corollary 2.12.** Suppose that a point  $y_o \in C$  satisfies the equation  $\langle x - y_o, y_o \rangle = 0$ and the inequality  $\langle x - y_o, y \rangle \leq 0$  for any  $y \in C$ . Then

$$\langle oldsymbol{x} - oldsymbol{y}_o, oldsymbol{y} - oldsymbol{y}_o, oldsymbol{y}_o, oldsymbol{y}_o 
angle \ \leq \ 0$$

for arbitrary  $m{y} \in m{C}$ . Consequently, according to Theorem 2.11 (b),  $m{y}_o = \Pi m{x}$ .

On the other hand, since C is a convex cone, for any number  $t \ge 0$  the point  $t \prod x$  belongs to C, too, whence

$$0 = \left. \frac{d}{dt} \right|_{t=1} \|\boldsymbol{x} - t \, \Pi \boldsymbol{x}\|^2 = -2 \langle \boldsymbol{x}, \Pi \boldsymbol{x} \rangle + 2 \|\Pi \boldsymbol{x}\|^2 = -2 \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \Pi \boldsymbol{x} \rangle.$$

Furthermore, for arbitrary  $y \in C$ , the sum  $\Pi x + y$  is an element of C, too, whence Theorem 2.11 (b) implies that

$$0 \geq \langle oldsymbol{x} - \Pi oldsymbol{x}, (\Pi oldsymbol{x} + oldsymbol{y}) - \Pi oldsymbol{x} 
angle \; = \; \langle oldsymbol{x} - \Pi oldsymbol{x}, oldsymbol{y} 
angle 
angle$$

For any  $x \in V$  and  $\lambda \ge 0$ , the point  $\lambda \Pi x \in C$  satisfies the (in)equalities

$$egin{aligned} &\langle \lambda m{x} - \lambda \Pi m{x}, \lambda \Pi m{x} 
angle &= \lambda^2 \langle m{x} - \Pi m{x}, \Pi m{x} 
angle &= 0, \ &\langle \lambda m{x} - \lambda \Pi m{x}, m{y} 
angle &= \lambda \langle m{x} - \Pi m{x}, m{y} 
angle &\leq 0 & ext{for all } m{y} \in m{C}. \end{aligned}$$

Hence  $\Pi(\lambda \boldsymbol{x}) = \lambda \Pi \boldsymbol{x}$ .

**Proof of Corollary 2.13.** Suppose that a point  $y_o \in C$  satisfies the equation  $\langle x - y_o, z \rangle = 0$  for arbitrary  $z \in C$ . Then  $\langle x - y_o, y - y_o \rangle = 0$  for arbitrary  $y \in C$ , and Theorem 2.11 (b) shows that  $y_o = \Pi x$ .

On the other hand,  $z := \Pi x \pm y \in C$  for arbitrary  $y \in C$ , whence Theorem 2.11 (b) implies that

$$0 \geq \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \boldsymbol{z} - \Pi \boldsymbol{x} \rangle = \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \pm \boldsymbol{y} \rangle = \pm \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \boldsymbol{y} \rangle$$

Consequently,  $\langle \boldsymbol{x} - \Pi \boldsymbol{x}, \boldsymbol{y} \rangle = 0$  for arbitrary  $\boldsymbol{y} \in \boldsymbol{C}$ .

For  $x_1, x_2 \in V$  and  $\lambda_1, \lambda_2 \in \mathbb{R}$ ,

$$\langle (\lambda_1 \boldsymbol{x}_1 + \lambda_2 \boldsymbol{x}_2) - (\lambda_1 \Pi \boldsymbol{x}_1 + \lambda_2 \Pi \boldsymbol{x}_2), \boldsymbol{y} \rangle = \lambda_1 \langle \boldsymbol{x}_1 - \Pi \boldsymbol{x}_1, \boldsymbol{y} \rangle + \lambda_2 \langle \boldsymbol{x}_2 - \Pi \boldsymbol{x}_2, \boldsymbol{y} \rangle = 0.$$

Hence  $\Pi(\lambda_1 \boldsymbol{x}_1 + \lambda \boldsymbol{x}_2) = \lambda_1 \Pi \boldsymbol{x}_1 + \lambda_2 \Pi \boldsymbol{x}_2$ . In other words,  $\Pi$  is a linear mapping from  $\boldsymbol{V}$  onto  $\boldsymbol{C}$ .

**Exercise 2.14** (Distance functions). Let A be a nonvoid subset of a metric space  $(\mathcal{X}, d)$ . For  $y \in \mathcal{X}$  let

$$d(y,A) := \inf_{x \in A} d(x,y).$$

Show that

- (a) d(y, A) = 0 if and only if y is in the closure  $\overline{A}$  of A,
- **(b)**  $d(y, A) = d(y, \overline{A}),$
- (c)  $|d(y,A) d(z,A)| \le d(y,z)$  for arbitrary  $y, z \in \mathcal{X}$ .

**Exercise 2.15.** Let  $V = \mathbb{R}^4$  and  $C := \{x \in \mathbb{R}^4 : x_1 \le x_2 \le x_3 \le x_4\}$ . Show that C is a closed convex cone. Guess the projection  $\Pi x$  of the vector  $x := (3, 1, 5, 3)^\top$  onto C. Verify your guess by checking the conditions in Corollary 2.12.

**Exercise 2.16** (Riesz'<sup>2</sup> representation theorem in Hilbert spaces). By means of Corollary 2.13 one can easily prove a key result in functional analysis: Let  $L : V \to \mathbb{R}$  be a continuous linear function. Then there exists a unique vector  $v \in V$  such that

$$L(\boldsymbol{x}) = \langle \boldsymbol{x}, \boldsymbol{v} \rangle$$
 for all  $\boldsymbol{x} \in \boldsymbol{V}$ .

Prove this result along the following lines:

(a) Show that  $C := \{x \in V : L(x) = 0\}$  is a closed linear subspace of V.

(b) If  $L \neq 0$ , there exists a  $x_1 \in V$  such that  $L(x_1) = 1$ . Now consider the vector  $v := \|x_1 - \Pi x_1\|^{-2} (x_1 - \Pi x_1)$  which satisfies  $v \perp C$  by Corollary 2.13.

(b.1) Verify that  $\langle \boldsymbol{x}_1, \boldsymbol{v} \rangle = 1$ .

(b.2) Show that for any  $x \in V$ , the vector  $x - L(x)x_1$  belongs to C. Deduce from the latter fact that  $L(x) = \langle x, v \rangle$ .

(c) Verify uniqueness of v.

<sup>&</sup>lt;sup>2</sup>Frigyes Riesz (1880–1956): Hungarian mathematician who made fundamental contributions to functional analysis.

### 2.2.2 Support functions and half spaces

Starting from Theorem 2.11 one can deduce various properties of convex sets. At first we show that a closed concex subset of V is always an intersection of closed halfspaces. A set  $H \subset V$  is called a closed halfspace if

$$oldsymbol{H} \;=\; ig\{oldsymbol{x} \in oldsymbol{V} : \langleoldsymbol{x},oldsymbol{v}
angle \leq rig\}$$

for some  $v \in V \setminus \{0\}$  and  $r \in \mathbb{R}$ . One can easily verify that this set H is closed and convex, see also Exercise 2.6. Its boundary is the hyperplane

$$ig\{oldsymbol{x}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{v}
angle=rig\}=rac{r}{\|oldsymbol{v}\|^2}oldsymbol{v}+oldsymbol{v}^\perp$$

with the orthogonal complement

$$\boldsymbol{v}^{\perp} = \left\{ \boldsymbol{y} \in \boldsymbol{V} : \langle \boldsymbol{y}, \boldsymbol{v} 
angle = 0 
ight\}$$

of v, a closed linear subspace of V.

**Definition 2.17** (Support function). The support function of a nonvoid set  $M \subset V$  is defined as

$$h_{\boldsymbol{M}}: \boldsymbol{V} \to (-\infty, \infty], \quad h_{\boldsymbol{M}}(\boldsymbol{v}) := \sup_{\boldsymbol{x} \in \boldsymbol{M}} \langle \boldsymbol{x}, \boldsymbol{v} \rangle.$$

Exercise 2.18 (Support functions of convex hulls). Show that

$$h_{\boldsymbol{M}} \equiv h_{\overline{\boldsymbol{M}}}$$
 and  $h_{\boldsymbol{M}} \equiv h_{\operatorname{conv}(\boldsymbol{M})},$ 

where  $\overline{M}$  denotes the closure of M. In particular, the support function of M coincides with the support function of  $\overline{\text{conv}(M)}$ , which is the smallest closed and convex set containing M, see also Exercise 2.8.

Recall that C is a nonvoid closed and convex subset of V and that B is the closed unit ball in V, so  $\partial B$  is the unit sphere  $\{u \in V : ||u|| = 1\}$ .

**Corollary 2.19** (Support functions and halfspaces). The support function  $h = h_C$  of C has the following properties:

$$oldsymbol{C} \;=\; igcap_{oldsymbol{v}\inoldsymbol{V}}ig\{oldsymbol{x}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{v}
angle\leq h(oldsymbol{v})ig\}\;=\;igcap_{oldsymbol{u}\in\partialoldsymbol{B}}ig\{oldsymbol{x}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{u}
angle\leq h(oldsymbol{u})ig\}.$$

For  $x \in V \setminus C$ ,

$$d(oldsymbol{x},oldsymbol{C}) \;=\; \max_{oldsymbol{u}\in\partialoldsymbol{B}}ig(\langleoldsymbol{x},oldsymbol{u}
angle - h(oldsymbol{u})ig).$$

Note that the set  $\{x \in V : \langle x, v \rangle \leq h(v)\}$  is a closed halfspace in V if  $v \neq 0$  and  $h(v) < \infty$ . Otherwise it is the full space V.

**Proof of Corollary 2.19.** Obviously, the set C is contained in any set  $\{x \in V : \langle x, v \rangle \leq h(v)\}$ , so

$$oldsymbol{C} \ \subset \ igcap_{oldsymbol{v}\inoldsymbol{V}}ig\{oldsymbol{x}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{v}
angle\leq h(oldsymbol{v})ig\} \ \subset \ igcap_{oldsymbol{u}\in\partialoldsymbol{B}}ig\{oldsymbol{x}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{u}
angle\leq h(oldsymbol{u})ig\}.$$

To show that the latter set coincides with C, it suffices to verify the particular representation of d(x, C) for  $x \in V \setminus C$ . To this end we consider the metric projection  $\Pi x$  of x onto C. For any unit vector u in V, the definition of h(u) and the Cauchy–Schwarz inequality imply that

$$\langle \boldsymbol{x}, \boldsymbol{u} \rangle - h(\boldsymbol{u}) \leq \langle \boldsymbol{x}, \boldsymbol{u} \rangle - \langle \Pi \boldsymbol{x}, \boldsymbol{u} \rangle = \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \boldsymbol{u} \rangle \leq \| \boldsymbol{x} - \Pi \boldsymbol{x} \| = d(\boldsymbol{x}, \boldsymbol{C}).$$

Specifically let  $u = ||x - \Pi x||^{-1} (x - \Pi x)$ . Then it follows from Theorem 2.11 (b) that

$$\langle \boldsymbol{y} - \Pi \boldsymbol{x}, \boldsymbol{u} \rangle = \| \boldsymbol{x} - \Pi \boldsymbol{x} \|^{-1} \langle \boldsymbol{y} - \Pi \boldsymbol{x}, \boldsymbol{x} - \Pi \boldsymbol{x} \rangle \le 0 \quad \text{for all } \boldsymbol{y} \in \boldsymbol{C}$$

In other words,  $\langle \boldsymbol{y}, \boldsymbol{u} \rangle \leq \langle \Pi \boldsymbol{x}, \boldsymbol{u} \rangle$  for all  $\boldsymbol{y} \in \boldsymbol{C}$ , whence

$$h(\boldsymbol{u}) = \sup_{\boldsymbol{y} \in \boldsymbol{C}} \langle \boldsymbol{y}, \boldsymbol{u} \rangle = \langle \Pi \boldsymbol{x}, \boldsymbol{u} \rangle,$$

and

$$\langle \boldsymbol{x}, \boldsymbol{u} \rangle = \langle \Pi \boldsymbol{x}, \boldsymbol{u} \rangle + \langle \boldsymbol{x} - \Pi \boldsymbol{x}, \boldsymbol{u} \rangle = h(\boldsymbol{u}) + \| \boldsymbol{x} - \Pi \boldsymbol{x} \|.$$

**Exercise 2.20** (Support functions and interior points). Show that for any  $x \in V$  and  $\delta > 0$  the following two statements are equivalent:

(i) The closed ball x + δB = {x + δv : v ∈ B} with center x and radius δ is contained in C;
(ii) h(u) ≥ ⟨x, u⟩ + δ for any unit vector u in V.

### 2.2.3 Separating hyperplanes

The next consequence of Theorem 2.11 is that for any point on the boundary of C there exists a tangent hyperplane, at least in case of V being finite-dimensional.

**Corollary 2.21** (Tangent hyperplanes). Let  $\dim(V) < \infty$ . For each  $x \in \partial C$  there exists a unit vector  $u \in V$  such that  $\langle x, u \rangle \ge \langle y, u \rangle$  for all  $y \in C$ .

**Proof of Corollary 2.21.** Let  $(x_n)_{n=1}^\infty$  be a sequence in  $V \setminus C$  with limit x. We define

$$u_n := \|x_n - \Pi x_n\|^{-1} (x_n - \Pi x_n).$$

Then  $(\boldsymbol{u}_n)_{n=1}^{\infty}$  is a sequence of unit vectors such that

$$\langle oldsymbol{y} - oldsymbol{x}_n, oldsymbol{u}_n 
angle = \langle oldsymbol{y} - \Pi oldsymbol{x}_n, oldsymbol{u}_n 
angle - \|oldsymbol{x}_n - \Pi oldsymbol{x}_n\| \ \leq \ 0 \ \ \ ext{for all } oldsymbol{y} \in oldsymbol{C},$$

see Theorem 2.11 (b). Since the unit sphere of a finite-dimensional Hilbert space is compact, we may assume without loss of generality that  $(u_n)_{n=1}^{\infty}$  converges to some unit vector u. But then

$$\langle oldsymbol{y},oldsymbol{u}
angle - \langle oldsymbol{x},oldsymbol{u}
angle \,=\, \langle oldsymbol{y}-oldsymbol{x},oldsymbol{u}
angle \,=\, \lim_{n o\infty} \langle oldsymbol{y}-oldsymbol{x}_n,oldsymbol{u}_n
angle \,\leq\, 0$$

for any  $y \in C$ .

**Theorem 2.22** (Separating hyperplanes). Let  $D_1$  and  $D_1$  be nonvoid convex subsets of V such that  $D_1 \cap D_2 = \emptyset$ .

(a) If  $D_1$  is compact and  $D_2$  is closed, then there exists a unit vector  $u \in V$  such that

(b) If  $\dim(V) < \infty$ , then there exists a unit vector  $u \in V$  such that

**Proof of Theorem 2.22.** We start with part (a). As shown in Exercise 2.14, the mapping  $x \mapsto d(x, D_2)$  is continuous. Hence compactness of  $D_1$  implies existence of a point  $a \in D_1$  such that  $d(a, D_2) \leq d(x, D_2)$  for all  $x \in D_1$ . If we denote the metric projection of a onto  $D_2$  with b, then

$$0 < \| \boldsymbol{b} - \boldsymbol{a} \| = \inf_{\boldsymbol{x} \in \boldsymbol{D}_1, \, \boldsymbol{y} \in \boldsymbol{D}_2} \, \| \boldsymbol{y} - \boldsymbol{x} \|.$$

In particular, a is the metric projection of b onto  $D_1$ . Hence it follows from Theorem 2.11 (b) applied to  $C = D_1, D_2$  that the unit vector

$$u := \|b - a\|^{-1}(b - a)$$

satisfies the following inequalities:

$$egin{array}{lll} \langle m{x}-m{a},m{u}
angle &\leq 0 & ext{for all }m{x}\inm{D}_1, \ \langle m{y}-m{b},m{u}
angle &\geq 0 & ext{for all }m{y}\inm{D}_2. \end{array}$$

This implies that

$$\sup_{oldsymbol{x}\inoldsymbol{D}_1}ig\langle oldsymbol{x},oldsymbol{u}ig
angle \ \leq \langleoldsymbol{a},oldsymbol{u}ig
angle, \ \sup_{oldsymbol{y}\inoldsymbol{D}_2}ig\langle oldsymbol{y},oldsymbol{u}ig
angle \ \geq \langleoldsymbol{b},oldsymbol{u}ig
angle,$$

and

$$\langle \boldsymbol{b}, \boldsymbol{u} 
angle - \langle \boldsymbol{a}, \boldsymbol{u} 
angle \ = \ \langle \boldsymbol{b} - \boldsymbol{a}, \boldsymbol{u} 
angle \ = \ \| \boldsymbol{b} - \boldsymbol{a} \| \ > \ 0.$$

Now to part (b). Since V is separable, which is a consequence of  $\dim(V) < \infty$ , for j = 1, 2 there exists a countable dense subset  $\{z_{j1}, z_{j2}, z_{j3}, \ldots\}$  of  $D_j$ ; see also Remark 2.23. It follows from Exercise 2.7 and convexity of  $D_j$  that  $D_{jn} := \operatorname{conv}\{z_{j1}, \ldots, z_{jn}\}$  is a compact subset of  $D_j$ . Moreover, this construction implies that  $D_{j1} \subset D_{j2} \subset D_{j3} \subset \cdots$  with

$$\lim_{n\to\infty} d(\boldsymbol{x},\boldsymbol{D}_{jn}) = 0 \quad \text{for any } \boldsymbol{x}\in\boldsymbol{D}_j.$$

According to part (a), for any index n there exists a unit vector  $u_n \in V$  such that

$$\inf_{oldsymbol{y}_n\inoldsymbol{D}_{2n}}ig\langleoldsymbol{y}_n,oldsymbol{u}_nig
angle\,>\,\sup_{oldsymbol{x}_n\inoldsymbol{D}_{1n}}ig\langleoldsymbol{x}_n,oldsymbol{u}_nig
angle.$$

Since dim $(V) < \infty$ , the unit sphere of V is compact. Hence we may replace  $((D_{1n}, D_{2n}, u_n))_n$ with a subsequence, if necessary, such that  $(u_n)_n$  converges to a unit vector u. But then for arbitrary  $x \in D_1$  and  $y \in D_2$  the metric projections  $x_n := \prod_{D_{1n}} x$  and  $y_n := \prod_{D_{2n}} y$  satisfy  $\lim_{n\to\infty} x_n = x$  and  $\lim_{n\to\infty} y_n = y$ , so the properties of  $(u_n)_n$  imply that

$$\langle \boldsymbol{y}, \boldsymbol{u} \rangle - \langle \boldsymbol{x}, \boldsymbol{u} \rangle = \lim_{n \to \infty} (\langle \boldsymbol{y}_n, \boldsymbol{u}_n \rangle - \langle \boldsymbol{x}_n, \boldsymbol{u}_n \rangle) \geq 0.$$

Consequently,  $\inf_{y \in D_2} \langle y, u \rangle$  is greater than or equal to  $\sup_{x \in D_1} \langle x, u \rangle$ .

**Remark 2.23** (Separable metric spaces). A metric space  $(\mathcal{X}, d)$  is called *separable* if there exists a countable set  $\mathcal{X}_o \subset \mathcal{X}$  which is dense in  $\mathcal{X}$ . The latter property means that for arbitrary  $x \in \mathcal{X}$ and  $\epsilon > 0$  there exists a  $x_o \in \mathcal{X}_o$  with  $d(x, x_o) \leq \epsilon$ . (A standard example is  $\mathcal{X} = \mathbb{R}^d$  with the standard Euclidean distance and its countable subset  $\mathcal{X}_o = \mathbb{Q}^d$ .)

For any nonvoid subset  $\mathcal{Y}$  of  $\mathcal{X}$ , the metric space  $(\mathcal{Y}, d)$  is separable, too. To see this, we 'choose' for each  $x_o \in \mathcal{X}_o$  and  $k \in \mathbb{N}$  a point  $y(x_o, k) \in \mathcal{Y}$  such that  $d(x_o, y(x_o, k)) \leq d(x_o, \mathcal{Y}) + 1/k$ . Then  $\mathcal{Y}_o := \{y(x_o, k) : x_o \in \mathcal{X}_o, k \in \mathbb{N}\}$  is a countable dense subset of  $\mathcal{Y}$ . Indeed, consider an arbitrary point  $y \in \mathcal{Y}$ . For  $\epsilon > 0$ , there exists a point  $x_o \in \mathcal{X}_o$  such that  $d(y, x_o) \leq \epsilon/3$ . Then for any integer  $k \geq 3/\epsilon$ , the point  $y(x_o, k) \in \mathcal{Y}_o$  satisfies

$$d(y, y(x_o, k)) \leq d(y, x_o) + d(x_o, y(x_o, k))$$
  
$$\leq d(y, x_o) + d(x_o, \mathcal{Y}) + 1/k$$
  
$$\leq 2d(y, x_o) + 1/k$$
  
$$\leq \epsilon.$$

#### 2.2.4 Polar sets and cones

For a subset M of V its polar set is defined as

$$oldsymbol{M}^* \ := \ ig\{oldsymbol{y} \in oldsymbol{V} : \langle oldsymbol{x}, oldsymbol{y} 
angle \leq 1 \ ext{ for all } oldsymbol{x} \in oldsymbol{M} ig\}.$$

Obviously,  $M^*$  contains the point 0, and writing

$$oldsymbol{M}^* \ = \ igcap_{oldsymbol{x}\inoldsymbol{M}}ig\{oldsymbol{y}\inoldsymbol{V}:\langleoldsymbol{x},oldsymbol{y}
angle\leq 1ig\}$$

shows that  $M^*$  is closed and convex. With the support function of M one may also write

$$\boldsymbol{M}^* = \{ \boldsymbol{y} \in \boldsymbol{V} : h_{\boldsymbol{M}}(\boldsymbol{y}) \leq 1 \}.$$

If the closed and convex set C contains the point 0, its polar set  $C^*$  has remarkable properties:

Theorem 2.24 (Polar sets). Suppose that C contains 0. Then

$$(\boldsymbol{C}^*)^* = \boldsymbol{C}$$

The set *C* is bounded if and only if 0 is an interior point of  $C^*$ . The set  $C^*$  is bounded if and only if 0 is an interior point of *C*.

$$oldsymbol{C} \ \subset \ ig\{oldsymbol{x} \in oldsymbol{V} : \langle oldsymbol{x}, oldsymbol{y} 
angle \leq 1 \ ext{ for all } oldsymbol{y} \in oldsymbol{C}^*ig\} \ = \ (oldsymbol{C}^*)^*.$$

It remains to show that for any  $x_1 \in V \setminus C$  there exists a  $y \in C^*$  such that  $\langle x_1, y \rangle > 1$ . To this end let  $x_0 := \prod_C x_1$ . Then  $v := x_1 - x_0 \neq 0$ , and by Theorem 2.11 (b),  $\langle x - x_0, v \rangle \leq 0$  for all  $x \in C$ . Consequently,

$$\sup_{oldsymbol{x}\inoldsymbol{C}} \langle oldsymbol{x},oldsymbol{v}
angle \ = \ \langle oldsymbol{x}_0,oldsymbol{v}
angle \ = \ \langle oldsymbol{x}_1,oldsymbol{v}
angle - \|oldsymbol{v}\|^2 \ < \ \langle oldsymbol{x}_1,oldsymbol{v}
angle.$$

Since  $\mathbf{0} \in \mathbf{C}$ , we know that

$$0 \leq \sup_{\boldsymbol{x} \in \boldsymbol{C}} \langle \boldsymbol{x}, \boldsymbol{v} 
angle \; = \; \langle \boldsymbol{x}_0, \boldsymbol{v} 
angle \; < \; \langle \boldsymbol{x}_1, \boldsymbol{v} 
angle.$$

In case of  $\langle \boldsymbol{x}_0, \boldsymbol{v} \rangle = 0$ , any point  $\boldsymbol{y} = t\boldsymbol{v}$  with  $t \ge 0$  belongs to  $\boldsymbol{C}^*$ , and  $\langle \boldsymbol{x}_1, \boldsymbol{y} \rangle = t \langle \boldsymbol{x}_1, \boldsymbol{v} \rangle > 1$ for sufficiently large t > 0. In case of  $\langle \boldsymbol{x}_0, \boldsymbol{v} \rangle > 0$  the point  $\boldsymbol{y} := \langle \boldsymbol{x}_0, \boldsymbol{v} \rangle^{-1} \boldsymbol{v}$  belongs to  $\boldsymbol{C}^*$  with  $\langle \boldsymbol{x}_1, \boldsymbol{y} \rangle > 1 = \langle \boldsymbol{x}_0, \boldsymbol{y} \rangle = \sup_{\boldsymbol{x} \in \boldsymbol{C}} \langle \boldsymbol{x}, \boldsymbol{y} \rangle$ .

Suppose that C is bounded, that means,  $C \subset rB$  for some r > 0 with the closed unit ball B of V. This implies that

$$\boldsymbol{C}^* \supset (r\boldsymbol{B})^* = r^{-1}\boldsymbol{B},$$

see Exercise 2.25. Thus 0 is an interior point of  $C^*$ . On the other hand, suppose that 0 is an interior point of C, that means,  $rB \subset C$  for some r > 0. Then

$$\boldsymbol{C}^* \subset (r\boldsymbol{B})^* = r^{-1}\boldsymbol{B},$$

whence  $C^*$  is bounded.

Since  $(C^*)^* = C$ , the previous two statements remain valid if we interchange C and  $C^*$ .

Exercise 2.25 (Polar sets of balls centered at 0). Show that for any r > 0 and the closed unit ball B of V,

$$(r\boldsymbol{B})^* = r^{-1}\boldsymbol{B}.$$

**Remark 2.26** (Polar cones). Suppose that C is a closed convex cone. Then one can show that

$$C^* = \{ \boldsymbol{y} \in \boldsymbol{V} : \langle \boldsymbol{x}, \boldsymbol{y} \rangle \leq 0 \}.$$

The latter set is a closed convex cone, too, the so-called *polar cone* of C.

Figure 2.1 indicates a closed convex cone  $C \subset \mathbb{R}^2$  and its polar cone  $C^*$ .

### **2.3 Extremal Points**

In this section we consider again an arbitrary real linear space V, and C is always a nonvoid convex subset of V. Recall the notion of convex hulls: For a nonvoid set  $M \subset V$  we defined its convex hull conv(M), the smallest convex set containing M. Now the question is whether for a given convex set C there is a *smallest* set  $M \subset V$  such that C = conv(M). If such a set Mexists, it will certainly contain all "extremal points" of C:


Figure 2.1: A closed convex cone C and its polar cone  $C^*$ .

**Definition 2.27** (Extremal points). A point  $x \in C$  is called extremal point of C (or extremal in C) if it satisfies the following condition: If  $x = (1 - \lambda)y + \lambda z$  with  $y, z \in C$  and  $0 < \lambda < 1$ , then y = z.

The set of all extremal points of C is denoted with extr(C). Any point  $x \in C \setminus extr(C)$  can be written as  $x = (1 - \lambda)y + \lambda z$  with two *different* points  $y, z \in C$  and some  $\lambda \in (0, 1)$ .

**Remark 2.28** (Simplified criterion for extremal points). A point  $x \in C$  is extremal in C if and only if it satisfies the following condition: If  $x = 2^{-1}(y + z)$  with  $y, z \in C$ , then y = z.

Proof: It is clear that any point  $x \in \text{extr}(C)$  satisfies the latter condition. On the other hand, suppose that  $x \in C$  satisfies the latter condition and equals  $(1 - \lambda)y + \lambda z$  with  $y, z \in C$  and  $\lambda \in (0, 1)$ . Without loss of generality let  $\lambda \leq 1/2$ , because otherwise we could interchange y and z. Then  $x = 2^{-1}(y + z')$  with  $z' = (1 - 2\lambda)y + 2\lambda z \in C$ . But then y = z', which is equivalent to y = z.

**Exercise 2.29.** Let x be an element of the convex set  $C \subset V$ . Show that  $x \in extr(C)$  if and only if  $C \setminus \{x\}$  is convex.

**Exercise 2.30**. Let  $C = \operatorname{conv}(M)$  for some set  $M \subset V$ . Show that

$$\operatorname{extr}(\boldsymbol{C}) \subset \boldsymbol{M}.$$

More precisely, show that

$$\operatorname{extr}(oldsymbol{C}) \ = \ ig\{oldsymbol{x}\in oldsymbol{M}:oldsymbol{x}
ot\in\operatorname{conv}(oldsymbol{M}\setminus\{oldsymbol{x}\})ig\}.$$

**Remark 2.31**. Suppose that  $(V, \|\cdot\|)$  is a normed space. Then

$$\operatorname{extr}(\boldsymbol{C}) \subset \partial \boldsymbol{C}.$$

Proof: If x is an interior point of C, then for some  $\epsilon > 0$ , the closed ball with center x and radius  $\epsilon$  is contained in C. Hence, for any  $w \in V$  with  $||w|| = \epsilon$ , the points y := x - w and z := x + w are different points in C such that  $x = 2^{-1}(y + z)$ .

**Exercise 2.32.** Let  $(V, \langle \cdot, \cdot \rangle, \| \cdot \|)$  be a Hilbert space. Show that any unit vector is an extremal point of the closed unit ball  $B = \{x \in V : ||x|| \le 1\}$ . Precisely,

$$\operatorname{extr}(\boldsymbol{B}) = \partial \boldsymbol{B}.$$

For compact and convex subsets of  $\mathbb{R}^d$  there is a simple representation theorem which answers our initial question:

**Theorem 2.33** (Minkowski). Let K be a compact and convex subset of  $\mathbb{R}^d$ . Then extr(K) is nonempty, and

$$\boldsymbol{K} = \operatorname{conv}(\operatorname{extr}(\boldsymbol{K})).$$

In Functional Analysis there exist extensions of this theorem to normed and topological vector spaces, e.g. the Krein–Milman theorem or Choquet's theorem. But for our purposes the finite-dimensional variant will be sufficient.

**Proof of Theorem 2.33.** In case of d = 1 the statement is easily verified: Here  $\mathbf{K} = [a, b]$  with real numbers  $a \leq b$ , and one can show easily that  $\operatorname{extr}(\mathbf{K}) = \{a, b\}$  and  $\mathbf{K} = \operatorname{conv}(\{a, b\})$ .

Suppose that for some d > 1 the assertion is correct for compact, convex subsets of  $\mathbb{R}^{d-1}$ . Now we have to show that an arbitrary point  $x \in K$  may be represented as a convex combination of points in extr(K).

Suppose first that x is in the interior of K. Fixing an arbitrary unit vector u, we may move from x in direction  $\pm u$  until we hit boundary points y = x - su and z = x + tu of K, where s, t > 0. But then  $x = (1 - \lambda)y + \lambda z$  with  $\lambda := s/(s + t) \in (0, 1)$ . Thus it suffices to consider boundary points x of K.

Let x be on the boundary of K. As shown in Corollary 2.21, there exists a unit vector u such that

$$r_o := \langle oldsymbol{x}, oldsymbol{u} 
angle \ = \ \max_{oldsymbol{y} \in oldsymbol{K}} \langle oldsymbol{y}, oldsymbol{u} 
angle.$$

In particular,  $\boldsymbol{x}$  belongs to the compact, convex set

$$oldsymbol{K}_o \ := \ ig\{oldsymbol{y} \in oldsymbol{K} : \langleoldsymbol{y},oldsymbol{u}
angle = r_oig\}.$$

But this set is contained in the hyperplane

$$oldsymbol{H}_o \ := \ ig\{oldsymbol{y} \in \mathbb{R}^d : \langle oldsymbol{y}, oldsymbol{u} 
angle = r_oig\} \ = \ ig\{r_ooldsymbol{u} + oldsymbol{z} : oldsymbol{z} \in oldsymbol{u}^otig\}$$

which is geometrically equivalent to  $\mathbb{R}^{d-1}$ . Consequently, by our induction hypothesis, x is a convex combination of extremal points of  $K_o$ . Now the assertion follows from the fact that any extremal point  $x_o$  of  $K_o$  is automatically an extremal point of K. For if  $x_o = 2^{-1}(y + z)$  with  $y, z \in K$ , then

$$r_o = \langle \boldsymbol{x}_o, \boldsymbol{u} 
angle = 2^{-1} \Big( \underbrace{\langle \boldsymbol{y}, \boldsymbol{u} 
angle}_{\leq r_o} + \underbrace{\langle \boldsymbol{z}, \boldsymbol{u} 
angle}_{\leq r_o} \Big),$$

whence  $\langle \boldsymbol{y}, \boldsymbol{u} \rangle = \langle \boldsymbol{z}, \boldsymbol{u} \rangle = r_o$ . Consequently,  $\boldsymbol{y}$  and  $\boldsymbol{z}$  belong to  $\boldsymbol{K}_o$ , so  $\boldsymbol{x}_o \in \text{extr}(\boldsymbol{K}_o)$  implies that  $\boldsymbol{y} = \boldsymbol{z} = \boldsymbol{x}_o$ .

**Exercise 2.34.** Prove or falsify the following claim for d = 1, 2, 3: If  $K \subset \mathbb{R}^d$  is compact and convex, then extr(K) is compact, too.

**Exercise 2.35.** Let  $C \subset [a, b]^d$  for real numbers a < b. Show that

$$oldsymbol{C} \cap \{a,b\}^d \ \subset \ \mathrm{extr}(oldsymbol{C}).$$

Then deduce from Exercise 2.30 that

$$\operatorname{extr}(\boldsymbol{C}) = \boldsymbol{C} \cap \{a, b\}^d \text{ if } \boldsymbol{C} = \operatorname{conv}(\boldsymbol{C} \cap \{a, b\}^d).$$

**Example 2.36**. Consider the following subsets of  $\mathbb{R}^d$ :

$$C_1 := \left\{ \boldsymbol{x} \in [0, \infty)^d : \sum_{i=1}^d x_i \le 1 \right\}, C_2 := \left\{ \boldsymbol{x} \in [0, 1]^d : x_1 \le x_2 \le \dots \le x_d \right\}.$$

Both sets are the intersection of finitely many closed halfspaces in  $\mathbb{R}^d$ , so they are closed and convex. (Precisely, both sets are the intersection of d + 1 closed halfspaces.) Moreover, both are contained in  $[0, 1]^d$ , so they are bounded and thus compact.

Now we show that for j = 1, 2,

$$\operatorname{extr}(\boldsymbol{C}_j) = \boldsymbol{C}_j \cap \{0, 1\}^d.$$

By Exercise 2.35 this is true, provided that  $C_i$  is the convex hull of  $C_i \cap \{0,1\}^d$ .

The set  $C_1 \cap \{0,1\}^d$  consists of the standard basis vectors  $e_1, \ldots, e_d$  and **0**. Any  $x \in C_1$  can be written as

$$x = \sum_{i=1}^{n} x_i e_i = \sum_{i=1}^{d} x_i e_i + x_{d+1} 0$$

with the nonnegative weights  $x_1, \ldots, x_d$  and  $x_{d+1} := 1 - \sum_{i=1}^d x_i$  summing to one. Thus  $C_1$  is indeed the convex hull of  $\{e_1, \ldots, e_d, \mathbf{0}\}$ .

Similarly, the set  $C_2 \cap \{0,1\}^d$  consists of  $b_1, b_2, \ldots, b_{d+1}$  with  $b_j := (1_{[i \ge j]})_{i=1}^d$ , i.e.  $b_{d+1} = 0$ . Then any  $x \in C_2$  may be written as

$$oldsymbol{x} \ = \ \sum_{j=1}^{d+1} \lambda_j oldsymbol{b}_j$$

with the nonnegative weights  $\lambda_1 := x_1$ ,  $\lambda_j := x_j - x_{j-1}$  for  $2 \le j \le d$  and  $\lambda_{d+1} := 1 - x_d$  summing to one.

The set  $C_2$  in the previous example may be seen as a special case of the set  $C \cap [0, 1]^d$  in the next exercise.

**Exercise 2.37.** Let  $\mathcal{O}$  be a nonempty set of index pairs (i, j), where  $i, j \in \{1, \ldots, d\}$  and  $i \neq j$ . Then define

$$C := \{ \boldsymbol{x} \in \mathbb{R}^d : x_i \le x_j \text{ whenever } (i, j) \in \mathcal{O} \}.$$

(a) Show that C is a closed convex cone.

(b) Show that

$$oldsymbol{C}\cap [0,1]^d = \operatorname{conv}(oldsymbol{E}_+) \quad ext{with} \quad oldsymbol{E}_+ \, := \, oldsymbol{C}\cap \{0,1\}^d.$$

(Hence, by Exercise 2.35, the elements of  $E_+$  are the extremal points of  $C \cap [0, 1]^d$ .)

(c) Show that

$$C = \operatorname{span}(\mathbf{1}_d) + C_+$$
 with  $C_+ := C \cap [0,\infty)^d$ .

(d) Show that

$$C_+ = \operatorname{cone}(E_+)$$
 and  $C = \operatorname{cone}(E_+ \cup \{-\mathbf{1}_d\})$ 

## 2.4 Convex Polyhedra

In this section we consider the space  $V = \mathbb{R}^d$  with its standard inner product. In the last section we encountered already various examples of a special type of subsets of  $\mathbb{R}^d$ :

**Definition 2.38** (Convex polyhedron). A nonvoid subset of  $\mathbb{R}^d$  is called a *convex polyhedron* if it is the intersection of finitely many closed halfspaces.

Since any closed halfspace is convex, a convex polyhedron is a closed and convex set. The next theorem provides some further facts which are important in various settings.

Theorem 2.39 (Properties of convex polyhedra). Consider a convex polyhedron

$$oldsymbol{C} \ = \ igcap_{i=1}^m ig\{ oldsymbol{x} \in \mathbb{R}^d : \langle oldsymbol{x}, oldsymbol{v}_i 
angle \leq r_i ig\}$$

with  $m \in \mathbb{N}$ ,  $v_i \in \mathbb{R}^d \setminus \{\mathbf{0}\}$  and  $r_i \in \mathbb{R}$ .

(a) The following three properties are equivalent:

(2.2)  $\max_{1 \le i \le m} \langle \boldsymbol{w}, \boldsymbol{v}_i \rangle > 0 \quad \text{for all } \boldsymbol{w} \in \mathbb{R}^d \setminus \{\mathbf{0}\};$ 

(2.3) **0** is an interior point of  $conv(v_1, \ldots, v_m)$ .

(b) A point  $x \in C$  is extremal in C if and only if

(2.4) 
$$\operatorname{span}\left\{\boldsymbol{v}_{i}: 1 \leq i \leq m, \langle \boldsymbol{x}, \boldsymbol{v}_{i} \rangle = r_{i}\right\} = \mathbb{R}^{d}.$$

In particular, extr(C) comprises at most  $\binom{m}{d}$  different points.

**Proof of Theorem 2.39.** We start with the equivalence of (2.1) and (2.2). Suppose that property (2.2) is violated. That means, for some  $w \in \mathbb{R}^d \setminus \{0\}$ ,

$$\langle \boldsymbol{w}, \boldsymbol{v}_i \rangle \leq 0 \quad \text{for } i \in \{1, \dots, m\}.$$

For a fixed point  $x_o \in C$  and arbitrary  $t \ge 0$ , the point  $x_o + tw$  lies in C as well, because  $\langle x_o + tw, v_i \rangle \le r_i + t \langle w, v_i \rangle \le r_i$  for i = 1, ..., m. Consequently, C is unbounded.

On the other hand, suppose that C is unbounded. Then there exists a sequence  $(x_n)_n$  in C with  $0 < ||x_n|| \to \infty$  as  $n \to \infty$ . We may replace this sequence with a subsequence, if necessary, such that  $\lim_{n\to\infty} ||x_n||^{-1}x_n = w$  for some unit vector  $w \in \mathbb{R}^d$ . But this vector satisfies the condition

$$\langle \boldsymbol{w}, \boldsymbol{v}_i 
angle = \lim_{n o \infty} \| \boldsymbol{x}_n \|^{-1} \underbrace{\langle \boldsymbol{x}_n, \boldsymbol{v}_i 
angle}_{\leq r_i} \leq \lim_{n o \infty} \| \boldsymbol{x}_n \|^{-1} r_i = 0$$

for arbitrary  $i \leq \{1, \ldots, m\}$ . Hence property (2.2) is violated as well.

Now we prove equivalence of (2.2) and (2.3). Let  $D = conv(v_1, \ldots, v_m)$ . As shown in Exercise 2.18,

$$h_{oldsymbol{D}}(oldsymbol{w}) := \sup_{oldsymbol{x} \in oldsymbol{D}} \langle oldsymbol{x}, oldsymbol{w} 
angle \ = \max_{1 \leq i \leq m} \langle oldsymbol{v}_i, oldsymbol{w} 
angle.$$

Consequently, (2.2) is equivalent to

$$h_{\boldsymbol{D}}(\boldsymbol{u}) > 0$$
 for any  $\boldsymbol{u} \in \partial \boldsymbol{B}$ .

Since  $h_D$  is continuous and  $\partial B$  is compact, this is even equivalent to

$$\min_{\boldsymbol{u}\in\partial\boldsymbol{B}}h_{\boldsymbol{D}}(\boldsymbol{u}) > 0.$$

But it follows from Exercise 2.20, that for any  $\delta > 0$ , the set  $\mathbf{0} + \delta \mathbf{B} = \delta \mathbf{B}$  is contained in  $\mathbf{D}$  if and only if  $h_{\mathbf{D}} \ge \delta$  on  $\partial \mathbf{B}$ . Thus (2.2) is equivalent to  $\mathbf{0}$  being an interior point of  $\mathbf{D}$ , which is (2.3).

Now we verify part (b). Suppose first that  $x \in C$  satisfies (2.4). To show that  $x \in \text{extr}(C)$ , let  $x = 2^{-1}(y + z)$  with  $y, z \in C$ . Now we have to show that y = z. Indeed,

$$\langle oldsymbol{x},oldsymbol{v}_i
angle \ = \ 2^{-1} \underbrace{\langleoldsymbol{y},oldsymbol{v}_i
angle}_{\leq r_i} + 2^{-1} \underbrace{\langleoldsymbol{z},oldsymbol{v}_i
angle}_{\leq r_i} < \ r_i \quad ext{if } \minig(\langleoldsymbol{y},oldsymbol{v}_i
angle,\langleoldsymbol{z},oldsymbol{v}_i
angle ig) \ < \ r_i,$$

whence  $\langle y - z, v_i \rangle = 0$  for all indices *i* such that  $\langle x, v_i \rangle = r_i$ . By assumption (2.4), these vectors  $v_i$  span the full space, so y - z = 0.

Secondly, suppose that  $x \in C$  violates condition (2.4). Then there exists a unit vector  $w \in \mathbb{R}^d$  such that  $\langle w, v_i \rangle = 0$  whenever  $\langle x, v_i \rangle = r_i$ . For  $\epsilon > 0$ ,

$$\langle \boldsymbol{x} \pm \epsilon \boldsymbol{w}, \boldsymbol{v}_i \rangle = \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle \pm \epsilon \langle \boldsymbol{w}, \boldsymbol{v}_i \rangle \leq r_i \quad \text{for all } i \in \{1, 2, \dots, m\},$$

provided that

$$\epsilon \leq \min \Big\{ rac{r_i - \langle \boldsymbol{x}, \boldsymbol{v}_i 
angle}{\left| \langle \boldsymbol{w}, \boldsymbol{v}_i 
angle \right|} : 1 \leq i \leq m, \langle \boldsymbol{x}, \boldsymbol{v}_i 
angle < r_i \Big\}.$$

In the latter case,  $x = 2^{-1}y + 2^{-1}z$  with  $y := x - \epsilon w \in C$  and  $z := x + \epsilon w \in C$ . Thus x is not an extremal point of C.

The upper bound for the cardinality of extr(C) may be verified as follows: If m < d, the set C has no extremal points at all. Otherwise one could determine all extremal points of C as follows (in principle): For arbitrary indices  $1 \le i(1) < \cdots < i(d) \le m$  check whether the vectors  $v_{i(1)}$ ,  $\ldots$ ,  $v_{i(d)}$  are linearly independent. If yes, let x be the unique vector satisfying  $\langle x, v_{i(j)} \rangle = r_{i(j)}$  for  $1 \le j \le d$ . Then check whether  $\langle x, v_i \rangle \le r_i$  for all remaining indices i. If yes, x is an extremal point of C. Since there are  $\binom{m}{d}$  possibilities for chosing  $i(1), \ldots, i(d)$ , there are at most  $\binom{m}{d}$  extremal points in C.

For compact convex polyhedra, there is an alternative simple representation:

**Theorem 2.40** (Compact convex polyhedra). A set  $C \subset \mathbb{R}^d$  is a compact convex polyhedron if and only if it is the convex hull of a finite subset of  $\mathbb{R}^d$ .

Before proving this theorem, we need a simple auxiliary result about convex sets.

**Lemma 2.41**. Let  $C = \operatorname{conv}(M)$  for some set  $M \subset \mathbb{R}^d$  with at least two elements. For any fixed  $x_o \in C$ , the following two properties of M are equivalent:

(i) The set C has nonempty interior.

(ii) The vectors  $x - x_o$ ,  $x \in M$ , span the whole space  $\mathbb{R}^d$ .

**Proof of Lemma 2.41.** Note first that  $C \subset x_o + \operatorname{span}(M - x_o)$ . Indeed, any point  $x \in C$  can be written as  $x = \sum_{i=1}^n \lambda_i x_i$  with  $n \in \mathbb{N}$  positive numbers  $\lambda_1, \ldots, \lambda_n$  summing to one and points  $x_1, \ldots, x_n \in M$ . But this may be rewritten as  $x = x_o + \sum_{i=1}^n \lambda_i (x_i - x_o)$ , a point in the affine space  $x_o + \operatorname{span}(M - x_o)$ .

Suppose that (ii) is violated, that means, span $(M - x_o)$  has dimension d' < d. Then there exists a unit vector  $u \in \mathbb{R}^d$  which is perpendicular to span $(M - x_o)$ , in other words, span $(M - x_o) \subset$  $u^{\perp}$ . In particular,  $C \subset x_o + \text{span}(M - x_o)$  is contained in the hyperplane  $x_o + u^{\perp}$ . Thus, Chas no interior points, whence (i) is violated, too.

Suppose that (ii) is satisfied. Then there exist points  $x_1, \ldots, x_d \in M$  such that  $x_1 - x_o, \ldots, x_d - x_o$  are linearly independent. In other words, the matrix  $A := [x_1 - x_o, \ldots, x_d - x_o] \in \mathbb{R}^{d \times d}$  is nonsingular, and  $\phi(\lambda) := x_o + A\lambda$  defines a homeomorphism  $\phi : \mathbb{R}^d \to \mathbb{R}^d$ . In particular, C contains the nonempty open set  $\phi(U)$ , where  $U := \{\lambda \in (0, \infty)^d : \sum_{i=1}^d \lambda_i < 1\}$ . Hence, (i) is satisfied as well.

**Proof of Theorem 2.40.** Theorem 2.33 implies that a compact convex polyhedron C is equal to  $\operatorname{conv}(\operatorname{extr}(C))$ . According to Theorem 2.39 (b), the set  $\operatorname{extr}(C)$  is finite. Thus C is the convex hull of a finite subset of  $\mathbb{R}^d$ .

On the other hand, suppose that  $C = \operatorname{conv}(x_1, \ldots, x_m)$  with  $m \in \mathbb{N}$  and pairwise different points  $x_i \in \mathbb{R}^d$ . In case of m = 1,  $C = \{x_1\}$  is easily seen to be the intersection of 2d closed halfspaces, so we only consider the case  $m \ge 2$ .

Suppose first that C has nonvoid interior. Without loss of generality let 0 be an interior point of C. Then it follows from Theorem 2.24 that the polar set  $C^* = \{y \in \mathbb{R}^d : h_C(y) \leq 1\}$  is also a compact and convex set with interior point 0. But since C is the convex hull of  $x_1, x_2, \ldots, x_m$ , it follows from Exercise 2.18 that

$$oldsymbol{C}^* \ = \ \Big\{oldsymbol{y} \in \mathbb{R}^d: \max_{1 \leq i \leq m} \langle oldsymbol{x}_i, oldsymbol{y} 
angle \leq 1 \Big\} \ = \ igcap_{i=1}^m ig\{oldsymbol{y} \in \mathbb{R}^d: \langle oldsymbol{x}_i, oldsymbol{y} 
angle \leq 1 ig\},$$

i.e.  $C^*$  is a compact convex polyhedron. According to the first part,  $C^*$  is the convex hull of a finite subset of  $\mathbb{R}^d$ . But then we may apply the previous considerations with  $C^*$  in place of C and  $(C^*)^* = C$  in place of  $C^*$  to show that C is a convex polyhedron, too.

In general, without loss of generality let  $\mathbf{0} \in C$ , and define

$$\boldsymbol{W} := \operatorname{span}\{\boldsymbol{x}_1,\ldots,\boldsymbol{x}_m\}.$$

Considering C as a subset of W, it has nonempty interior, see Lemma 2.41. Thus C may be represented as

$$oldsymbol{C} \;=\; igcap_{i=1}^n ig\{ oldsymbol{x} \in oldsymbol{W}: \langle oldsymbol{x}, oldsymbol{v}_i 
angle \leq r_i ig\}$$

with  $n \in \mathbb{N}$  nonzero vectors  $v_1, \ldots, v_n \in W$  and real numbers  $r_1, \ldots, r_n$ . But if  $b_1, \ldots, b_k$  is an orthonormal basis of  $W^{\perp}$ , then we may rewrite C as

$$oldsymbol{C} \ = igcap_{i=1}^n ig\{ oldsymbol{x} \in \mathbb{R}^d : \langle oldsymbol{x}, oldsymbol{v}_i 
angle \leq r_i ig\} \ \cap \ igcap_{j=1}^k ig\{ oldsymbol{x} \in \mathbb{R}^d : \langle oldsymbol{x}, oldsymbol{b}_j 
angle \leq 0 ig\} \cap igcap_{j=1}^k ig\{ oldsymbol{x} \in \mathbb{R}^d : \langle oldsymbol{x}, -oldsymbol{b}_j 
angle \leq 0 ig\},$$

so *C* is a convex polygon.

**Exercise 2.42**. Let M be an arbitrary nonvoid subset of a real vector space V. Show that for any fixed  $x_o \in \text{conv}(M)$  the following three sets are identical:

$$egin{aligned} m{W}_1 &:= ext{span}\{m{x} - m{x}_o : m{x} \in m{M}\}, \ m{W}_2 &:= ext{span}\{m{x} - m{y} : m{x}, m{y} \in m{M}\}, \ m{W}_3 &:= ext{span}igg\{\sum_{i=1}^n \lambda_i m{x}_i : n \in \mathbb{N}; m{x}_1, \dots, m{x}_n \in m{M}; \lambda_1, \dots, \lambda_n \in \mathbb{R} ext{ with } \sum_{i=1}^n \lambda_i = 0igg\}. \end{aligned}$$

#### 2.4.1 Linear programs

In numerous applications one has to minimize a linear function  $f : \mathbb{R}^d \to \mathbb{R}$ ,

$$(2.5) f(\boldsymbol{x}) = \langle \boldsymbol{x}, \boldsymbol{c} \rangle$$

with a given nonzero vector  $c \in \mathbb{R}^d$ , under certain constraints on x. Precisely, let x be an element of

(2.6) 
$$\boldsymbol{C} := \left\{ \boldsymbol{x} \in [0,\infty)^d : \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b} \right\}$$

with a given matrix  $A = (a_{i,j})_{i \le k,j \le d} \in \mathbb{R}^{k \times d}$  of rank k < d and a given vector  $b \in \mathbb{R}^k$ . This set C is a convex polyhedron determined by d + 2k linear inequalities:

$$\langle \boldsymbol{x}, \pm \boldsymbol{a}_i \rangle \leq \pm b_i \quad \text{for } 1 \leq i \leq k,$$
  
 $\langle \boldsymbol{x}, -\boldsymbol{e}_j \rangle \leq 0 \quad \text{for } 1 \leq j \leq d,$ 

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_k]^{\top}$  with linearly independent vectors  $\mathbf{a}_1, \dots, \mathbf{a}_k \in \mathbb{R}^d$ , and  $\mathbf{e}_1, \dots, \mathbf{e}_d$  denotes the standard basis of  $\mathbb{R}^d$ . The following corollary to Theorem 2.39 provides additional information about C and its extremal points.

Corollary 2.43. Let C be the polyhedron defined in (2.6). This set is bounded if and only if

$$\left\{ oldsymbol{w} \in [0,\infty)^d : oldsymbol{A}oldsymbol{w} = oldsymbol{\{0\}} = oldsymbol{\{0\}}$$

A point  $x \in C$  is extremal in C if and only if

$$\operatorname{span}(\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_k\}\cup\{\boldsymbol{e}_i:x_i=0\}) = \mathbb{R}^d.$$

**Proof of Corollary 2.43.** Recall that C is a convex polyhedron determined by m = 2k + d linear inequalities with directions  $\pm a_1, \ldots, \pm a_k$  and  $-e_1, \ldots, -e_d$ . Suppose that  $w \in \mathbb{R}^d$  satisfies  $\langle w, \pm a_j \rangle \leq 0$  for  $1 \leq j \leq k$  and  $\langle w, -e_i \rangle \leq 0$  for  $1 \leq i \leq d$ . This is equivalent to saying that Aw = 0 and  $w \in [0, \infty)^d$ . Hence, by Theorem 2.39 (a), C is bounded if and only if w = 0 is the only vector with these properties.

The characterization of  $x \in \text{extr}(C)$  is an immediate consequence of Theorem 2.39 (b).

In case of *C* being compact, we know that C = conv(extr(C)), and

$$\inf_{\boldsymbol{x}\in\boldsymbol{C}}f(\boldsymbol{x}) = \min_{\boldsymbol{x}\in \operatorname{extr}(\boldsymbol{C})}f(\boldsymbol{x}).$$

This equation can be verified easily. It also follows from Exercise 2.18 and the fact that

$$\inf_{\boldsymbol{x}\in\boldsymbol{C}}f(\boldsymbol{x})\ =\ -h_{\boldsymbol{C}}(-\boldsymbol{c}).$$

Even if C is unbounded, either the minimum is attained at an extremal point of C or f is unbounded from below:

**Theorem 2.44** (Solutions of a linear program). Let f and C be given by (2.5) and (2.6), respectively. If there exists a vector  $w \in [0, \infty)^d$  with Aw = 0 and f(w) < 0, then

$$\inf_{\boldsymbol{x}\in\boldsymbol{C}} f(\boldsymbol{x}) = -\infty.$$

Otherwise,

$$\inf_{\boldsymbol{x} \in \boldsymbol{C}} f(\boldsymbol{x}) = \min_{\boldsymbol{x} \in \operatorname{extr}(\boldsymbol{C})} f(\boldsymbol{x}) \in \mathbb{R}.$$

**Proof of Theorem 2.44.** Suppose first that Aw = 0 and f(w) < 0 for some  $w \in [0, \infty)^d$ . For any  $x \in C$  and arbitrary  $t \ge 0$ , the point x + tw belongs to C as well, and

$$f(\boldsymbol{x} + t\boldsymbol{w}) = f(\boldsymbol{x}) + tf(\boldsymbol{w}) \rightarrow -\infty \text{ as } t \rightarrow \infty.$$

Thus

$$\inf_{\boldsymbol{y}\in\boldsymbol{C}}\,f(\boldsymbol{y})\,=\,-\infty$$

Now suppose that  $f(w) \ge 0$  whenever  $w \in [0, \infty)^d$  and Aw = 0. Suppose that  $x \in C$  is not an extremal point of C. That means,

$$\operatorname{span}(\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_k\}\cup\{\boldsymbol{e}_j:x_j=0\})\neq\mathbb{R}^d.$$

Hence there exists a unit vector w which is perpendicular to the latter space. In other words, w is a unit vector such that

$$Aw = 0$$
 and  $w_i = 0$  whenever  $x_i = 0$ .

Replacing  $\boldsymbol{w}$  with  $-\boldsymbol{w}$ , if necessary, we may assume without loss of generality that  $f(\boldsymbol{w}) \ge 0$  and  $w_j > 0$  for at least one index j. (Indeed, if  $f(\boldsymbol{w}) > 0$  but  $\boldsymbol{w} \in (-\infty, 0]^d$ , then  $\tilde{\boldsymbol{w}} := -\boldsymbol{w}$  would satisfy  $A\tilde{\boldsymbol{w}} = \mathbf{0}, \tilde{\boldsymbol{w}} \in [0, \infty)^d$  and  $f(\tilde{\boldsymbol{w}}) < 0$ .) But then

$$\boldsymbol{x}^{ ext{new}} := \boldsymbol{x} - t \boldsymbol{w} \quad ext{with} \quad t := \min\left\{rac{x_j}{w_j} : w_j > 0
ight\} > 0.$$

would be a new vector in C with  $f(\boldsymbol{x}^{\text{new}}) \leq f(\boldsymbol{x})$  and

$$\operatorname{span}(\{a_1,\ldots,a_k\}\cup\{e_j:x_j=0\}) \subseteq \operatorname{span}(\{a_1,\ldots,a_k\}\cup\{e_j:x_j^{\operatorname{new}}=0\})$$

Replacing x with  $x^{\text{new}}$  and repeating this step, if necessary, we reach eventually a point  $x \in \text{extr}(C)$  with the same or a smaller value of f(x) than the starting point. Since extr(C) is finite, this shows that

$$\inf_{\boldsymbol{x}\in\boldsymbol{C}} f(\boldsymbol{x}) = \min_{\boldsymbol{x}\in\operatorname{extr}(\boldsymbol{C})} f(\boldsymbol{x}) \in \mathbb{R}.$$

**Explicit formulae for extremal points of** C**.** As shown in Corollary 2.43, a point  $x \in C$  is extremal in C if and only if

$$\operatorname{span}(\{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_k\}\cup\{\boldsymbol{e}_j:x_j=0\}) = \mathbb{R}^d.$$

By means of basis completion one can show that this condition is equivalent to the following one: There exist indices  $1 \le j(1) < \cdots < j(d-k) \le d$  such that

$$\det([a_1,\ldots,a_k, e_{j(1)},\ldots,e_{j(d-k)}]) \neq 0 \text{ and } x_{j(1)} = \cdots = x_{j(d-k)} = 0.$$

If  $h(1) < \cdots < h(k)$  are the elements of  $\{1, \ldots, d\} \setminus \{j(1), \ldots, j(d-k)\}$ , then the former determinant equals  $\pm \det((a_{i,h(\ell)})_{i,\ell=1}^k))$ . Hence we arrive at the following conclusion:

A point  $x \in C$  is extremal in C if and only if there exist indizes h(1) < ... < h(k) in  $\{1, ..., d\}$  such that

$$\det\left((a_{i,h(\ell)})_{i,\ell=1}^k\right) \neq 0 \quad \text{and} \quad x_j = 0 \text{ if } j \notin \{h(1),\ldots,h(k)\}.$$

In this case,

$$(x_{h(\ell)})_{\ell=1}^k = ((a_{i,h(\ell)})_{i,\ell=1}^k)^{-1} \boldsymbol{b}.$$

The latter representation is at the heart of the so-called *simplex method*, invented by G.B. Dantzig in 1947. Roughly speaking, one starts with one extremal point  $x^{\mathcal{H}}$  of C, where  $\mathcal{H}$  stands for the index set  $\{h(1), \ldots, h(k)\}$ , and then one modifies  $\mathcal{H}$  step by step by exchanging one of its elements such that the value  $f(x^{\mathcal{H}})$  decreases.

**Example 2.45** (Nutrients of meals). Let  $F_1, F_2, \ldots, F_d$  be given food ingredients,  $d \ge 4$ . For  $1 \le i \le d$  let  $a_{1i}, a_{2i}, a_{3i} \ge 0$  be the weight fractions of proteins, carbohydrates and fat, respectively, contained in one weight unit of  $F_i$ . Suppose one is planning to compose a meal, given by a vector  $x \in [0, \infty)^d$ , where  $x_i$  specifies the amount (weight) of  $F_i$ . The goal is to have a meal containing given amounts  $b_1, b_2, b_3 > 0$  of proteins, carbohydrates and fat, respectively. The set of vectors x with these properties corresponds to the set C above, where k = 3. It follows from Corollary 2.43 that this set is bounded, unless  $a_{1i} = a_{2i} = a_{3i} = 0$  for some ingredient  $F_j$ .

Now one could think about finding a 'meal'  $x \in C$  with a maximal amount of one's favourite ingredient  $F_{j_*}$ , i.e. one would like to minimize  $f(x) := \langle x, -e_{j_*} \rangle$ . Alternatively, let  $c_j$  be the price of  $F_j$  per unit weight. Then  $f(x) = \langle x, c \rangle$  would be the overall price of the 'meal' x. One could also consider the volume  $c_j$  of one weight unit of  $F_j$  and thus try to minimize the overall volume of a meal. With  $c = \mathbf{1}_d := (1, ..., 1)^{\top}$  one would consider the overall weight of a meal.

Obviously, the extremal points of C may fail to be 'delicious' themselves; their main purpose is to determine the minimum of f over C. A data example will be presented in the lecture; see also the supplementary files.

**Exercise 2.46.** For a given program of study and a certain time interval, let  $p_i$  be the relative frequency of degrees with final grade  $x_i = 4 + (i - 1)/2$ ,  $1 \le i \le 5$ , among all degrees awarded. Thus  $p \in [0, 1]^5$  with  $\sum_{i=1}^5 p_i = 1$ . Suppose we only know that

$$\sum_{i=1}^{5} p_i x_i = \mu \in (4,6) \text{ and } \left(\sum_{i=1}^{5} p_i (x_i - \mu)^2\right)^{1/2} = \sigma > 0.$$

(a) Show that the set  $P(\mu, \sigma)$  of all vectors p with these properties is a compact polyhedron with at most 10 extremal points.

(b) Write a program which generates for given values  $\mu$  and  $\sigma$  a list of the extremal points of  $P(\mu, \sigma)$ .

(c) Determine with your program sharp bounds for  $p_j$ ,  $1 \le j \le 5$ , and  $\sum_{i=1}^k p_i$ ,  $1 \le k \le 4$ , in case of  $\mu = 5.5$  and  $\sigma = 0.5$ .

**Inequality constraints and slack variables.** Suppose we want to replace some of the constraints  $a_i^{\top} x = b_i$  with the weaker constraint  $a_i^{\top} x \leq b_i$ . After rearranging the rows of A, if necessary, let the constraint set C consist of all vectors  $x \in [0, \infty)^d$  such that

$$\mathbf{a}_i^\top \mathbf{x} \leq b_i \quad \text{if } 1 \leq i \leq k', \\ \mathbf{a}_i^\top \mathbf{x} = b_i \quad \text{if } k' < i \leq k.$$

Then we may introduce a vector  $\boldsymbol{y} \in [0,\infty)^{k'}$  of additional "slack variables" and write

$$b_i = \boldsymbol{a}_i^\top \boldsymbol{x} + y_i \quad \text{for } i = 1, \dots, k'$$

Thus we consider the augmented vector  $\tilde{\boldsymbol{x}} = [\boldsymbol{x}^{\top}, \boldsymbol{y}]^{\top} \in [0, \infty)^{d+k'}$  and want to minimize  $\tilde{f}(\tilde{\boldsymbol{x}}) := f(\boldsymbol{x})$  under the constraint that  $\tilde{\boldsymbol{A}}\tilde{\boldsymbol{x}} = \boldsymbol{b}$ , where

$$oldsymbol{A} \,:=\, [oldsymbol{A},oldsymbol{D}]$$
 with  $oldsymbol{D} \,:=\, (1_{[i=s]})_{i\leq k,s\leq k'}$ 

Hence by increasing the dimension d we may translate the weaker constraint set in a polyhedron of type (2.6).

#### 2.4.2 Doubly stochastic matrices and the Hoffmann–Wielandt inequalities

**Doubly stochastic matrices.** We consider a special example of a convex polyhedron. The matrix space  $\mathbb{R}^{d \times d}$  becomes a  $d^2$ -dimensional Euclidean space, if we define  $\langle \mathbf{A}, \mathbf{B} \rangle := \sum_{i,j=1}^{d} A_{ij} B_{ij} =$ trace $(\mathbf{A}^{\top} \mathbf{B})$ . The corresponding norm is the so-called Frobenius norm  $\|\cdot\|_F$ .

Now we consider the set  $\Theta$  of all matrices  $\boldsymbol{\theta} \in [0,\infty)^{d \times d}$  whose row and column sums are all equal to one,

$$\sum_{j=1}^{d} \theta_{ij} = 1 \text{ for } 1 \leq i \leq d \text{ and } \sum_{i=1}^{d} \theta_{ij} = 1 \text{ for } 1 \leq j \leq d.$$

This set is characterized by  $d^2 + 4d$  linear inequalities, and since it is contained in  $[0, 1]^d$ , it is compact. Hence our general results imply that  $extr(\Theta)$  is a finite set, and  $\Theta = conv(extr(\Theta))$ . The set of extremal points in  $\Theta$  may be characterized as follows:

**Theorem 2.47** (Birkhoff–von Neumann). The set  $extr(\Theta)$  is the set of all permutation matrices  $\theta$ , that means,

$$\boldsymbol{\theta} = \left( \mathbf{1}_{[j=\sigma(i)]} \right)_{i,j=1}^{d}$$

for some permutation  $\sigma$  of  $\{1, \ldots, d\}$ .

In other words,  $\boldsymbol{\theta} \in \mathbb{R}^{d \times d}$  is an extremal point of  $\Theta$  if and only if each row and each column contains precisely one entry 1 and d-1 entries 0.

**Proof of Theorem 2.47.** Instead of using Theorem 2.39 (b) we provide an elementary proof. Note first that  $\Theta \cap \{0, 1\}^{d \times d}$  is the set of all permutation matrices: The definition of  $\Theta$  entails that any matrix in  $\{0, 1\}^{d \times d}$  belongs to  $\Theta$  if and only if it has precisely on entry 1 and d - 1 entries 0 in each row and in each column, that means, it is a permutation matrix.

As shown in Exercise 2.35, any matrix in  $\theta \in \Theta \cap \{0, 1\}^{d \times d}$  is extremal in  $\Theta$ . Hence, it suffices to show that any matrix  $\theta$  in  $\Theta \setminus \{0, 1\}^{d \times d}$  is not extremal. Indeed, the set  $\mathcal{K} := \{(i, j) : 0 < \theta_{ij} < 1\}$  is nonempty. The definition of  $\Theta$  implies that with  $(i, j) \in \mathcal{K}$  also  $(i', j), (i, j') \in \mathcal{K}$  for some indices  $i' \neq i$  and  $j' \neq j$ . From this one can deduce that there is an even number  $k \ge 4$  of different pairs  $(i_1, j_1), \ldots, (i_k, j_k) \in \mathcal{K}$  such that

$$i_t = i_{t+1}$$
 for odd  $t \in \{1, \dots, k-1\}$ ,  
 $j_t = j_{t+1}$  for even  $t \in \{2, 3, \dots, k\}$  with  $j_{k+1} := j_1$ .



Figure 2.2: Two possible configurations of  $(i_1, j_1), \ldots, (i_k, j_k)$ .

Figure 2.2 illustrates the definition of this sequence  $(i_1, j_1), \ldots, (i_k, j_k)$ . Now we define the non-zero matrix  $\boldsymbol{\delta} \in \mathbb{R}^{d \times d}$  with components

$$\delta_{ij} := \begin{cases} (-1)^t & \text{if } (i,j) = (i_t, j_t), 1 \le t \le k, \\ 0 & \text{else.} \end{cases}$$

All row and column sums of this matrix  $\delta$  are equal to 0, and for sufficiently small  $\epsilon > 0$ ,

$$\boldsymbol{\theta} = 2^{-1} \big( (\boldsymbol{\theta} + \epsilon \boldsymbol{\delta}) + (\boldsymbol{\theta} - \epsilon \boldsymbol{\delta}) \big) \quad \text{and} \quad \boldsymbol{\theta} \pm \epsilon \boldsymbol{\delta} \in \Theta.$$

Thus  $\boldsymbol{\theta}$  is not an extremal point in  $\Theta$ .

**The Hoffmann–Wielandt inequalities.** The preceding considerations enable an elegant proof of two famous inequalities about eigenvalues of symmetric matrices.

**Theorem 2.48** (Hoffmann–Wielandt). Let A and B be symmetric matrices in  $\mathbb{R}^{d \times d}$  with eigenvalues  $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_d$  and  $\mu_1 \leq \mu_2 \leq \cdots \leq \mu_d$ , respectively. Then

$$\sum_{i=1}^{d} (\lambda_i - \mu_i)^2 \leq \sum_{i,j=1}^{d} (A_{ij} - B_{ij})^2 \leq \sum_{i=1}^{d} (\lambda_i - \mu_{d+1-i})^2.$$

Both inequalities are sharp. Just consider the diagonal matrices  $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_d)$  and  $\mathbf{B} = \text{diag}(\mu_1, \dots, \mu_d), \text{diag}(\mu_d, \dots, \mu_1)$ . The first inequality may be interpreted as a Lipschitz property of eigenvalues: The standard Euclidean distance between the two vectors  $\boldsymbol{\lambda}$  and  $\boldsymbol{\mu}$  of ordered eigenvalues is bounded from above by the Frobenius norm of  $\mathbf{A} - \mathbf{B}$ ,

$$\|\boldsymbol{\lambda} - \boldsymbol{\mu}\| \leq \|\boldsymbol{A} - \boldsymbol{B}\|_F.$$

**Proof of Theorem 2.48.** The spectral decomposition of symmetric matrices shows that there exist orthogonal matrices  $U, V \in \mathbb{R}^{d \times d}$  such that

$$oldsymbol{A} = oldsymbol{U} \operatorname{diag}(oldsymbol{\lambda}) oldsymbol{U}^ op$$
 and  $oldsymbol{B} = oldsymbol{V} \operatorname{diag}(oldsymbol{\mu}) oldsymbol{V}^ op$ 

The general rule trace(DE) = trace(ED) for arbitrary matrices  $D, E \in \mathbb{R}^{d \times d}$  implies that  $\|U^{\top}D\|_{F} = \|D\|_{F} = \|DV\|_{F}$ . In particular,

$$\begin{split} \|\boldsymbol{A} - \boldsymbol{B}\|_{F}^{2} &= \|\boldsymbol{U}^{\top} \boldsymbol{A} \boldsymbol{V} - \boldsymbol{U}^{\top} \boldsymbol{B} \boldsymbol{V}\|_{F}^{2} \\ &= \|\operatorname{diag}(\boldsymbol{\lambda}) \boldsymbol{U}^{\top} \boldsymbol{V} - \boldsymbol{U}^{\top} \boldsymbol{V} \operatorname{diag}(\boldsymbol{\mu})\|_{F}^{2} \\ &= \sum_{i,j=1}^{d} (\lambda_{i} - \mu_{j})^{2} (\boldsymbol{U}^{\top} \boldsymbol{V})_{ij}^{2} \\ &= \sum_{i,j=1}^{d} (\lambda_{i} - \mu_{j})^{2} \theta_{ij}, \end{split}$$

with  $\boldsymbol{\theta} := (\theta_{ij})_{i,j=1}^d$  given by

$$heta_{ij} \ := \ (oldsymbol{U}^ opoldsymbol{V})_{ij}^2$$

Hence for fixed vectors  $\lambda, \mu$ , the squared norm  $\|A - B\|_F^2$  is a linear function of the matrix  $\theta$ . Note that  $W := U^{\top}V$  is orthogonal, too, so

$$\sum_{j=1}^{d} \theta_{ij} = (\boldsymbol{W}\boldsymbol{W}^{\top})_{ii} = 1 \quad \text{for } 1 \le i \le d,$$
$$\sum_{i=1}^{d} \theta_{ij} = (\boldsymbol{W}^{\top}\boldsymbol{W})_{jj} = 1 \quad \text{for } 1 \le j \le d.$$

That means,  $\boldsymbol{\theta}$  belongs to the set  $\boldsymbol{\Theta}$  of all doubly stochastic matrices in  $\mathbb{R}^{d \times d}$ . As stated in Theorem 2.47,  $\boldsymbol{\Theta} = \operatorname{conv}(\boldsymbol{P})$  with  $\boldsymbol{P}$  the set of all permutation matrices in  $\mathbb{R}^{d \times d}$ . Thus

$$\min_{\sigma \in \mathcal{S}_d} \sum_{i=1}^d (\lambda_i - \mu_{\sigma(i)})^2 \leq \|\boldsymbol{A} - \boldsymbol{B}\|_F^2 \leq \max_{\sigma \in \mathcal{S}_d} \sum_{i=1}^d (\lambda_i - \mu_{\sigma(i)})^2$$

with the set  $S_d$  of all permutations of  $\{1, 2, ..., d\}$ . Now the assertion follows from Lemma 2.49 below.

**Lemma 2.49** (Hardy–Littlewood–Polya). For arbitrary real numbers  $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_d$  and  $\mu_1 \leq \mu_2 \leq \cdots \leq \mu_d$  and any permutation  $\sigma$  of  $\{1, 2, \ldots, d\}$ ,

$$\sum_{i=1}^{d} (\lambda_i - \mu_i)^2 \leq \sum_{i=1}^{d} (\lambda_i - \mu_{\sigma(i)})^2 \leq \sum_{i=1}^{d} (\lambda_i - \mu_{d+1-i})^2.$$

**Proof of Lemma 2.49.** This lemma is essentially a consequence of the fact that any permutation may be represented as a combination of finitely many transpositions (i.e. exchanges of pairs). For fixed indices  $1 \le j < k \le d$  we define a new permutation  $\tilde{\sigma}$  via

$$\tilde{\sigma}(i) := \begin{cases} \sigma(i) & \text{for } i \notin \{j, k\} \\ \sigma(k) & \text{for } i = j, \\ \sigma(j) & \text{for } i = k. \end{cases}$$

Then

$$\sum_{i=1}^{d} (\lambda_i - \mu_{\tilde{\sigma}(i)})^2 - \sum_{i=1}^{d} (\lambda_i - \mu_{\sigma(i)})^2$$
$$= (\lambda_j - \mu_{\sigma(k)})^2 + (\lambda_k - \mu_{\sigma(j)})^2 - (\lambda_j - \mu_{\sigma(j)})^2 - (\lambda_k - \mu_{\sigma(k)})^2$$
$$= 2(\lambda_k - \lambda_j)(\mu_{\sigma(k)} - \mu_{\sigma(j)}) \begin{cases} \leq 0 & \text{if } \sigma(k) < \sigma(j), \\ \geq 0 & \text{if } \sigma(k) > \sigma(j). \end{cases}$$

On the one hand, let j be the smallest index such that  $\sigma(j) \neq j$ . Then  $\sigma(j) > j$  and  $k := \sigma^{-1}(j) > j$  while  $\sigma(k) < \sigma(j)$ . Then we replace  $\sigma$  with  $\tilde{\sigma}$ , so  $\sigma(i) = i$  for  $i \leq j$ . Repeating this step as often as necessary, we end up with  $\sigma(i) = i$  for all i. In each step the sum of all squares  $(\lambda_i - \mu_{\sigma(i)})^2$  stays constant or decreases by (2.7). This yields the first asserted inequality.

On the other hand, let j be the smallest index such that  $\sigma(j) \neq d+1-j$ . Then  $\sigma(j) < d+1-j$ and  $k := \sigma^{-1}(d+1-j) > j$  while  $\sigma(k) > \sigma(j)$ . Then we replace  $\sigma$  with  $\tilde{\sigma}$ , so  $\sigma(i) = d+1-i$ for  $i \leq j$ . This step is repeated until finally  $\sigma(i) = d+1-i$  for all i. Since in each step the sum of all squares  $(\lambda_i - \mu_{\sigma(i)})^2$  stays constant or increases by (2.7), this yields the second asserted inequality.

### 2.4.3 Polyhedral cones

So far, we analyzed bounded convex polyhedra in quite some detail. Now we consider another special subfamily of convex polyhedra.

**Definition 2.50** (Polyhedral cone). A nonvoid subset  $\mathbb{R}^d$  is called a *polyhedral cone* if it is the intersection of finitely many closed halfspaces of the form  $\{x \in \mathbb{R}^d : \langle x, v \rangle \leq 0\}$  with  $v \in \mathbb{R}^d \setminus \{0\}$ .

Here is an analogue to Theorem 2.40 for polyhedral cones.

**Theorem 2.51** (Polyhedral cones). A set  $C \subset \mathbb{R}^d$  is a polyhedral cone if and only if it is the cone generated by a finite subset of  $\mathbb{R}^d$ .

This theorem can be deduced from Theorem 2.40 by means of a special decomposition of  $\mathbb{R}^d$  which is explained first. Let

(2.8) 
$$\boldsymbol{C} = \bigcap_{i=1}^{m} \{ \boldsymbol{x} \in \mathbb{R}^{d} : \langle \boldsymbol{x}, \boldsymbol{v}_{i} \rangle \leq 0 \}$$

with  $m \in \mathbb{N}$  and nonzero vectors  $\boldsymbol{v}_1, \ldots, \boldsymbol{v}_m \in \mathbb{R}^d$ . Now let

$$oldsymbol{V} := \operatorname{span}(oldsymbol{v}_1,\ldots,oldsymbol{v}_m) \quad ext{and} \quad oldsymbol{W} := oldsymbol{V}^\perp = \{oldsymbol{v}_1,\ldots,oldsymbol{v}_m\}^\perp.$$

Obviously, W is a subset of C. If  $x \in \mathbb{R}^d$  is written as  $x = x_V + x_W$  with  $x_V \in V$  and  $x_W \in W$ , then  $\langle x, v_i \rangle = \langle x_V, v_i \rangle$  for  $1 \le i \le m$ , whence  $x \in C$  if and only if  $x_V \in C$ . This shows that

$$C = C \cap V + W.$$

Note that  $C \cap V$  is a "pointed cone" in the sense that for  $x \in V \setminus \{0\}$ ,  $x \in C$  implies that  $-x \notin C$ .

**Proof of Theorem 2.51.** Suppose first that *C* is given by (2.8). With the decomposition  $C = C \cap V + W$  just mentioned, it suffices to show that  $C \cap V$  and *W* are cones generated by finite sets. Since  $\{0\} = \text{cone}(0)$ , it suffices to consider the cases  $C \cap V \neq \{0\}$  and  $W \neq \{0\}$ .

Let  $b_1, \ldots, b_q$  be a basis of W, and let  $b_0 := -\sum_{i=1}^q b_i$ . Any point  $x \in W$  may be written as  $x = \sum_{i=1}^q \mu_i b_i$  with suitable numbers  $\mu_1, \ldots, \mu_q \in \mathbb{R}$ . But for any  $\lambda_0 \ge 0$ ,

$$\boldsymbol{x} = \lambda_0 \boldsymbol{b}_0 + \sum_{i=1}^q (\mu_i + \lambda_0) \boldsymbol{b}_i,$$

and for sufficiently large  $\lambda_0$ , all factors  $\mu_i + \lambda_0$ ,  $1 \leq i \leq q$ , are positive. Hence,  $\mathbf{W} = \operatorname{cone}(\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_q)$ .

Suppose that  $C \cap V \neq \{0\}$ . With  $b := -\sum_{i=1}^{m} v_i$ , every nonzero vector  $x \in C \cap V$  satisfies  $\langle x, b \rangle = 1$ . Indeed, if  $x \in C$ , then  $\langle x, b \rangle = -\sum_{i=1}^{m} \langle x, v_i \rangle \ge 0$  with equality if and only if  $\langle x, v_i \rangle = 0$  for  $1 \le i \le m$ , whence  $x \in V^{\perp} = W$ . Consequently,  $C \cap V$  may be written as

$$\boldsymbol{C} \cap \boldsymbol{V} = \{\lambda \boldsymbol{x} : \lambda \geq 0, \boldsymbol{x} \in \boldsymbol{K}\}$$

with

$$oldsymbol{K}$$
 :=  $ig\{oldsymbol{x} \in oldsymbol{C} \cap oldsymbol{V} : \langle oldsymbol{x}, oldsymbol{b} 
angle = 1ig\}$ 

Note that K is a convex polyhedron, determined by  $m + 2\dim(W) + 2$  linear inequalities. If K is bounded, it follows from Theorem 2.40 that  $K = \operatorname{conv}(M)$  for some finite subset of  $C \cap V$ , whence  $C \cap V = \operatorname{cone}(M)$ . To verify boundedness of K, assume the contrary. Then there exists a sequence  $(x_n)_n$  in K such that  $||x_n|| \to \infty$  as  $n \to \infty$ . Without loss of generality we may assume that  $u_n := ||x_n||^{-1}x_n$  converges to a unit vector u as  $n \to \infty$ . Since  $C \cap V$  is a closed convex cone, all points  $u_n$  and the limit u belong to  $C \cap V$ . But this would imply the contradiction that

$$\langle \boldsymbol{u}, \boldsymbol{b} \rangle = \lim_{n \to \infty} \langle \boldsymbol{u}_n, \boldsymbol{b} \rangle = \lim_{n \to \infty} \|\boldsymbol{x}_n\|^{-1} = 0.$$

Now suppose that  $C = \operatorname{cone}(v_1, \ldots, v_m)$  for some  $m \in \mathbb{N}$  and nonzero vectors  $v_1, \ldots, v_m \in \mathbb{R}^d$ . As shown in Exercise 2.10, C is closed, whence  $C = (C^*)^*$  with  $M^*$  denoting the polar set of  $M \subset \mathbb{R}^d$ . But

$$oldsymbol{C}^* \ = \ ig\{oldsymbol{y} \in \mathbb{R}^d : \langleoldsymbol{y},oldsymbol{x}
angle \leq 0 ext{ for all }oldsymbol{x} \in oldsymbol{C}ig\} \ = \ ig\{oldsymbol{y} \in \mathbb{R}^d : \langleoldsymbol{y},oldsymbol{v}_i
angle \leq 0 ext{ for }1 \leq i \leq mig\},$$

which is a polyhedral cone. As shown before, there exists a finite subset M of  $\mathbb{R}^d$  such that  $C^* = \operatorname{cone}(M)$ . But then,

$$\boldsymbol{C} = \big\{ \boldsymbol{x} \in \mathbb{R}^d : \langle \boldsymbol{x}, \boldsymbol{y} \rangle \leq 0 \text{ for all } \boldsymbol{y} \in \operatorname{cone}(\boldsymbol{M}) \big\} = \big\{ \boldsymbol{x} \in \mathbb{R}^d : \langle \boldsymbol{x}, \boldsymbol{y} \rangle \leq 0 \text{ for all } \boldsymbol{y} \in \boldsymbol{M} \big\},$$

which is a polyhedral cone as well.

The polar cone of a polyhedral cone. It follows from the general theory that for any closed convex cone  $C \subset \mathbb{R}^d$  and its polar cone  $C^*$ , every point  $x \in \mathbb{R}^d$  can be written as  $x = \Pi x + \Pi^* x$ , where  $\Pi$  and  $\Pi^*$  denote the metric projections onto C and  $C^*$ , respectively. Moreover,  $\Pi x \perp \Pi^* x$ .

Now suppose that C is a polyhedral cone as in (2.8). Then,

$$\boldsymbol{C}^* = \operatorname{cone}(\boldsymbol{v}_1,\ldots,\boldsymbol{v}_m).$$

For  $\boldsymbol{x} \in \mathbb{R}^d$  let

$$J(\boldsymbol{x}) := \{i \in \{1, \dots, m\} : \langle \Pi \boldsymbol{x}, \boldsymbol{v}_i \rangle = 0\}.$$

Then

(2.9) 
$$\Pi \boldsymbol{x} \in \{\boldsymbol{v}_j : j \in J(\boldsymbol{x})\}^{\perp} \text{ and } \Pi^* \boldsymbol{x} \in \operatorname{cone}(\boldsymbol{v}_j : j \in J(\boldsymbol{x})),$$

where  $\operatorname{cone}(\emptyset) := \{\mathbf{0}\}$ . This implies that  $\Pi x$  is the orthogonal projection onto the linear space  $\{v_j : j \in J(x)\}^{\perp}$  and  $\Pi^* x$  is the orthogonal projection onto the linear space  $\operatorname{span}(v_j : j \in J(x))$ . The first part of (2.9) is just a consequence of the definition of J(x). As to the second part of (2.9), let  $\Pi^* x = \sum_{i=1}^m \lambda_i v_i$  with  $\lambda_1, \ldots, \lambda_m \geq 0$ . For any index j with  $\lambda_j > 0$ , the vector  $\Pi^* x - \lambda_j v_j$  belongs to  $C^*$  too, whence

$$0 \leq -\langle \Pi \boldsymbol{x}, \boldsymbol{v}_j \rangle = \lambda_j^{-1} \langle \Pi \boldsymbol{x}, \Pi^* \boldsymbol{x} - \lambda_j \boldsymbol{v}_j \rangle \leq 0,$$

whence  $j \in J(\boldsymbol{x})$ . Consequently,  $\Pi^* \boldsymbol{x} \in \operatorname{cone}(\boldsymbol{v}_j : j \in J(\boldsymbol{x}))$ .

## **Chapter 3**

# **Convex Functions**

## **3.1** Convex Functions on the Real Line

In Section 1.2 we encountered already convex functions. There, the functions f under consideration were at least differentiable, and convexity was just an additional feature. Now we review the definition of convexity of univariate functions and its implications in full generality and detail. Throughout this section let C be a nondegenerate interval in  $\mathbb{R}$ .

**Definition 3.1** (Convex and concave functions, I). (a) A function  $f : C \to (-\infty, \infty]$  is called *convex*, if for arbitrary points  $x, y \in C$  and  $\lambda \in (0, 1)$ ,

$$f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y).$$

A function  $f: C \to [-\infty, \infty)$  is called *concave*, if -f is convex.

(b) A function  $f : C \to \mathbb{R}$  is called *strictly convex*, if for arbitrary different points  $x, y \in C$  and  $\lambda \in (0, 1)$ ,

$$f((1-\lambda)x + \lambda y) < (1-\lambda)f(x) + \lambda f(y).$$

A function  $f: C \to \mathbb{R}$  is called *strictly concave*, if -f is strictly convex.

Note that for points x < y in C, the set  $\{(1 - \lambda)x + \lambda y : \lambda \in (0, 1)\}$  is the interval (x, y). Hence, if we define the *domain* of  $f : C \to (-\infty, \infty]$  to be the set

$$\operatorname{dom}(f) := \{ x \in C : f(x) < \infty \},\$$

then convexity of f implies that dom(f) is a subinterval of C. If we set  $f(x) := \infty$  for  $x \in \mathbb{R} \setminus C$ , then dom(f) remains unaltered, and convexity of f on C implies convexity of f on  $\mathbb{R}$ .

**Exercise 3.2** (Simplified criterion for convexity). Suppose that  $f : C \to \mathbb{R}$  is continuous. Show that f is convex (strictly convex) if and only if

$$f(2^{-1}(x+y)) \le (<) 2^{-1}(f(x)+f(y))$$

for arbitrary different points  $x, y \in C$ .

**Exercise 3.3** (A contraction principle). Let  $f : [-1,1] \to \mathbb{R}$  be convex. Show that for any  $x \in (-1,1)$ ,

$$f(x) + f(-x) \leq f(1) + f(-1).$$

Alternative characterizations. For three numbers x < y < z in C one can always write

$$y = (1 - \lambda)x + \lambda z$$
 with  $\lambda := \frac{y - x}{z - x} \in (0, 1)$  and  $1 - \lambda = \frac{z - y}{z - x}$ 

Hence, f is (strictly) convex on C if and only if for any choice of the points x < y < z in C,

$$f(y) \leq (<) \ (1-\lambda)f(x) + \lambda f(z).$$

Now we want to reinterpret such inequalities in terms of the slopes (difference ratios)

$$f'(r,s) := \frac{f(s) - f(r)}{s - r} \in [-\infty, \infty]$$

for arbitrary points  $r, s \in C$  such that  $r \neq s$  and at least one of them is in dom(f). Note that f'(r,s) = f'(s,r).

In case of  $f(y) < \infty$ , the inequality

(3.1) 
$$f(y) \leq (1-\lambda)f(x) + \lambda f(z)$$

is equivalent to

$$(3.2) f'(x,y) \leq f'(y,z).$$

Indeed,

$$(1-\lambda)f(x) + \lambda f(z) - f(y) = (1-\lambda)(f(x) - f(y)) + \lambda(f(z) - f(y))$$
  
= 
$$\frac{(z-y)(f(x) - f(y)) + (y-x)(f(z) - f(y))}{z-x}$$
  
= 
$$\frac{(y-x)(z-y)}{z-x}(f'(y,z) - f'(x,y)).$$

In case of  $f(x) < \infty$ , condition (3.1) is equivalent to

(3.3) 
$$f'(x,y) \leq f'(x,z).$$

Indeed, if  $f(z) = \infty$ , then both (3.1) and (3.3) are trivial. If  $f(z) < \infty$ , then both (3.1) and (3.3) imply that  $f(y) < \infty$ , and

$$f'(x,z) - f'(x,y) = \frac{1}{y-x} (\lambda(f(z) - f(x)) - f(y) + f(x))$$
  
=  $\frac{1}{y-x} ((1-\lambda)f(x) + \lambda f(z) - f(y)).$ 

Analogously one can show that in case of  $f(z) < \infty$ , condition (3.1) is equivalent to

(3.4) 
$$f'(x,z) \leq f'(y,z).$$

These considerations lead to the following characterizations of convexity.

**Lemma 3.4** (Characterizing convexity). For a function  $f : \mathbb{R} \to (-\infty, \infty]$  the following conditions are equivalent:

(a) f is convex;

(b) dom(f) is an interval, and for arbitrary points x < y < z in dom(f), inequality (3.2) holds true;

(c) for arbitrary points  $x \in dom(f)$ , the function

$$\mathbb{R} \setminus \{x\} \ni y \mapsto f'(x,y) \in [-\infty,\infty]$$

is isotonic.

With obvious modifications one can also verify the following characterization of strict convexity:

**Lemma 3.5** (Characterizing strict convexity). For a function  $f : C \to \mathbb{R}$ , the following two conditions are equivalent:

- (a) f is strictly convex;
- (b) for arbitrary points x < y < z in C, inequality (3.2) holds strictly;
- (c) for arbitrary points  $x \in C$ , the function

$$C \setminus \{x\} \ni y \mapsto f'(x,y) \in \mathbb{R}$$

is strictly isotonic.

**Smoothness and integral representations.** It follows from characterization (c) in Lemma 3.4 that any convex function f admits finite one-sided derivatives on the interior of dom(f), and the latter satisfy certain inequalities.

**Corollary 3.6** (One-sided derivatives). (a) Let  $f : \mathbb{R} \to (-\infty, \infty]$  be convex. Then for each  $x \in \text{dom}(f)$  the one-sided derivatives

$$f'(x-) := \lim_{y \to x, y < x} f'(x,y),$$
  
$$f'(x+) := \lim_{y \to x, y > x} f'(x,y)$$

exist in  $[-\infty, \infty]$ , and

 $f'(x-) \leq f'(x+).$ 

Moreover, for arbitrary  $x, y \in \text{dom}(f)$  with x < y:

$$f'(x+) \leq f'(x,y) \leq f'(y-).$$

In particular,  $f'(x \pm) \in \mathbb{R}$  for any interior point x of dom(f).

(b) Let  $f : C \to \mathbb{R}$  be strictly convex. Then for arbitrary  $x, y \in C$  with x < y,

$$f'(x+) < f'(x,y) < f'(y-).$$

**Remark 3.7.** For a convex function  $f : \mathbb{R} \to (-\infty, \infty]$ , Corollary 3.6 (a) has the following implications:

(i) The functions  $x \mapsto f'(x \pm)$  are monotonically increasing on dom(f) and real-valued on the interior of dom(f).

(ii) There are at most countably many points  $x \in \text{dom}(f)$  such that f'(x-) < f'(x+).

For if  $C_o$  is the set of all such points  $x \in \text{dom}(f)$ , then the intervals (f'(x-), f'(x+)),  $x \in C_o$ , are pairwise disjoint. Consequently, if we choose a rational number  $q(x) \in (f'(x-), f'(x+))$ , then  $C_o \ni x \mapsto q(x) \in \mathbb{Q}$  is injective.

(iii) For any nondegenerate interval  $[a, b] \subset \operatorname{dom}(f)$ ,

$$|f(y) - f(x)| \leq L|y - x|$$
 for  $x, y \in [a, b]$ ,

where  $L := \max\{f'(b-), -f'(a+)\}$ . In particular, f is Lipschitz-continuous on [a, b] whenever a and b are interior points of dom(f).

To see this, note that for  $a \le x < y \le b$ ,

$$f'(x,y) \begin{cases} \leq f'(x,b) \leq f'(b-), \\ \geq f'(a,y) \geq f'(a+). \end{cases}$$

Together with the mean value theorem, Lemmas 3.4, 3.5 and Corollary 3.6 imply the following facts for smooth functions:

**Lemma 3.8** (Convexity of differentiable functions). Let  $f : C \to \mathbb{R}$  be continuous, and let  $C^o$  be the interior of C.

(a) If f is differentiable on  $C^{\circ}$ , it is (strictly) convex if and only if f' is (strictly) isotonic on  $C^{\circ}$ .

(b) If f is twice differentiable on  $C^o$ , it is convex if and only if  $f'' \ge 0$  on  $C^o$ .

(c) If f is twice differentiable on  $C^o$ , it is strictly convex if and only if  $f'' \ge 0$  on  $C^o$  and  $\{x \in C^o : f''(x) > 0\}$  is dense in  $C^o$ .

**Proof of Lemma 3.8.** Suppose that  $f : C \to \mathbb{R}$  is (strictly) convex. Then by Corollary 3.6, f' is (strictly) isotonic on  $C^o$ . On the other hand, suppose that f' is (strictly) isotonic. According to the mean value theorem, for arbitrary points x < y < z in C, there exist points  $\xi_1 \in (x, y)$  and  $\xi_2 \in (y, z)$  such that

$$f'(x,y) = f'(\xi_1) \le (<) f'(\xi_2) = f'(y,z).$$

Thus it follows from Lemma 3.4 (b) (Lemma 3.5 (b)) that f is (strictly) convex. This proves part (a).

As to part (b), if f is twice differentiable, then f' is isotonic if and only if  $f'' \ge 0$ , so by part (a), f is convex if and only if  $f'' \ge 0$ .

As to part (c), by parts (a) and (b), we may assume that f' is isotonic and  $f'' \ge 0$ . Now we have to show that f' is strictly isotonic if and only if  $\{x \in C^o : f''(x) > 0\}$  is dense in  $C^o$ . If the latter condition is violated, then f'' = 0 on an interval  $[a, b] \subset C^o$  with a < b. But then, f' is constant on [a, b], so f' is not strictly increasing. If  $\{x \in C^o : f''(x) > 0\}$  is dense in  $C^o$ , then for arbitrary  $a, b \in C^o$  with a < b, there exists a point  $x \in (a, b)$  such that f''(x) > 0. But then for sufficiently small  $\epsilon > 0, x \pm \epsilon \in [a, b]$  and

$$f'(a) \leq f'(x-\epsilon) < f'(x+\epsilon) \leq f'(b).$$

This shows that f' is strictly increasing on  $C^o$ .

**Exercise 3.9.** Determine for each of the following functions  $f : \mathbb{R} \to \mathbb{R}$  maximal intervals C on which it is convex or concave<sup>1</sup>.

$$f(x) := \sqrt{1+x^2}; \qquad f(x) := \frac{x^2}{1+x^2}; \\ f(x) := \frac{x^2}{\sqrt{1+x^2}}; \qquad f(x) := \log(1+x^2); \\ f(x) := \log(1+e^x).$$

**Example 3.10**. For  $x \in \mathbb{R}$ , let f(x) := |x|. With the triangle inequality one can verify rigorously that f is convex with dom $(f) = \mathbb{R}$ . Moreover,

$$f'(x\pm) = \begin{cases} \operatorname{sign}(x) & \text{if } x \neq 0, \\ \pm 1 & \text{if } x = 0. \end{cases}$$

As a second example consider

$$f(x) \ := \ \begin{cases} -\sqrt{1-x^2} & \text{if } |x| \le 1, \\ \infty & \text{if } |x| > 1. \end{cases}$$

Here dom(f) = [-1, 1]. Moreover, f is convex on  $\mathbb{R}$  and strictly convex on [-1, 1], because  $f < 0 = f(\pm 1)$  on (-1, 1), and for  $x \in (-1, 1)$ ,

$$f'(x) = x(1-x^2)^{-1/2}, \quad f''(x) = (1-x^2)^{-3/2},$$

Moreover,

$$f'(x\pm) = \begin{cases} -\infty & \text{if } x = -1, \\ +\infty & \text{if } x = +1. \end{cases}$$

**Example 3.11** (An exotic convex function). Let  $Q = \{q_1, q_2, q_3, ...\}$  be a countable subset of  $\mathbb{R}$ , and let  $a_1, a_2, a_3, ... > 0$  such that  $\sum_{k=1}^{\infty} a_k(1 + |q_k|) < \infty$ . Then,

$$f(x) := \sum_{k=1}^{\infty} a_k |x - q_k|$$

defines a convex function  $f : \mathbb{R} \to \mathbb{R}$  such that f is differentiable at  $x \in \mathbb{R}$  if and only if  $x \in Q$ . Moreover, if  $Q \cap C$  is dense in some nondegenerate interval  $C \subset \mathbb{R}$ , then f is strictly convex on C. The verification of these properties is left to the reader as an exercise.

Our last finding about convex functions on the real line ist the fact that they are absolutely continuous on the interior of their domain:

<sup>1</sup>i.e. -f is convex

57

**Lemma 3.12** (Absolute continuity). Let  $f : \mathbb{R} \to (-\infty, \infty]$  be convex. For  $x \in \text{dom}(f)$  let f'(x) be an arbitrary number in [f'(x-), f'(x+)]. Then f' is isotonic on dom(f), and for arbitrary interior points a, b of dom(f),

$$f(b) - f(a) = \int_a^b f'(x) \, dx.$$

**Proof of Lemma 3.12.** Isotonicity of f' follows from Corollary 3.6. Without loss of generality let a < b. For arbitrary  $k \in \mathbb{N}$  and  $i \in \{0, ..., k\}$  we define  $\delta_k := (b - a)/k$  and  $t_{k,i} := a + i\delta_k$ . Then it follows from Corollary 3.6 and the construction of f' that

$$f(b) - f(a) = \sum_{i=0}^{k-1} (f(t_{k,i+1}) - f(t_{k,i})) = \sum_{i=0}^{k-1} f'(t_{k,i}, t_{k,i+1}) \delta_k \begin{cases} \geq \sum_{i=0}^{k-1} f'(t_{k,i}) \delta_k, \\ \leq \sum_{i=1}^{k} f'(t_{k,i}) \delta_k. \end{cases}$$

Now, isotonicity of f' implies that

$$\sum_{i=0}^{k-1} f'(t_{k,i})\delta_k = \sum_{i=1}^k f'(t_{k,i})\delta_k - (f'(b) - f'(a))\delta_k$$
$$\geq \sum_{i=1}^k \int_{t_{k,i-1}}^{t_{k,i}} f'(x) \, dx - (f'(b) - f'(a))\delta_k$$
$$= \int_a^b f'(x) \, dx - (f'(b) - f'(a))\delta_k$$

and

$$\sum_{j=1}^{k} f'(t_{k,j})\delta_k = \sum_{j=0}^{k-1} f'(t_{k,j})\delta_k + (f'(b) - f'(a))\delta_k$$
$$\leq \int_a^b f'(x) \, dx + (f'(b) - f'(a))\delta_k$$

As  $k \to \infty$ , it follows that  $f(b) - f(a) = \int_a^b f'(x) \, dx$ .

**Exercise 3.13.** Prove or disprove the following claim: A function  $f : \mathbb{R} \to \mathbb{R}$  with f(0) = 0 is convex if the function  $\mathbb{R} \setminus \{0\} \ni x \mapsto f(x)/x$  is isotonic.

**Exercise 3.14**. Let  $h : C \to \mathbb{R}$  be an isotonic (a strictly isotonic) function. For any fixed  $x_o \in C$  define  $f : C \to \mathbb{R}$  via

$$f(x) := \int_{x_o}^x h(t) \, dt.$$

Show that f is (strictly) convex with  $f'(x \pm) = h(x pm)$ .

**Exercise 3.15**. Consider functions  $g : \mathbb{R} \to \mathbb{R}$  and  $h : \mathbb{R} \to (-\infty, \infty]$ .

- (a) Show that  $h \circ g$  is convex, provided that g is convex and h is convex and isotonic.
- (b) Why is it essential in part (a) that h is isotonic?

- (c) Suppose that h is convex and antitonic. Which property of g would guarantee convexity of  $h \circ g$ ?
- (d) Find conditions on g and h which guarantee that  $h \circ g$  is strictly convex.

## **3.2** Convex Functions on Linear Spaces

Throughout this section let C be a nondegenerate convex subset of a real vector space V.

**Definition 3.16** (Convex functions, II). (a) A function  $f : \mathbf{C} \to (-\infty, \infty]$  is called *convex* if for arbitrary  $\mathbf{x}, \mathbf{y} \in \mathbf{C}$  and  $\lambda \in (0, 1)$ ,

$$f((1 - \lambda)\boldsymbol{x} + \lambda \boldsymbol{y}) \leq (1 - \lambda)f(\boldsymbol{x}) + \lambda f(\boldsymbol{y})$$

A function  $f : \mathbb{C} \to [-\infty, \infty)$  is called *concave* is -f is convex.

(b) A function  $f : C \to \mathbb{R}$  is called strictly convex if for arbitrary different points  $x, y \in C$  and  $\lambda \in (0, 1)$ ,

$$f((1-\lambda)\boldsymbol{x}+\lambda\boldsymbol{y}) < (1-\lambda)f(\boldsymbol{x})+\lambda f(\boldsymbol{y}).$$

A function  $f : C \to \mathbb{R}$  is called *strictly concave* is -f is strictly convex.

**Exercise 3.17** (Norms). Let  $(V, \|\cdot\|)$  a real normed vector space. Show that for any  $x_o \in V$ ,  $f(x) := \|x - x_o\|$  defines a convex function f on  $\mathbb{R}$ .

**Exercise 3.18** (Sublevel sets). Suppose that  $f : \mathbb{C} \to (-\infty, \infty]$  is convex. Show that all sublevel sets  $\{x \in \mathbb{C} : f(x) \le r\}$  and  $\{x \in \mathbb{C} : f(x) < r\}$  with  $r \in \mathbb{R}$  are convex. Does the latter property imply convexity?

Again we define

$$\operatorname{dom}(f) := \{ \boldsymbol{x} \in \boldsymbol{C} : f(\boldsymbol{x}) < \infty \}$$

for a function  $f : \mathbf{C} \to (-\infty, \infty]$ . If f is convex, then  $\operatorname{dom}(f)$  is a convex subset of  $\mathbf{C}$ . If we define  $f(\mathbf{x}) := \infty$  for  $\mathbf{x} \in \mathbf{V} \setminus \mathbf{C}$ , then convexity of f on  $\mathbf{C}$  is equivalent to convexity of f on  $\mathbf{V}$  while  $\operatorname{dom}(f)$  remains unaltered.

**Remark 3.19** (Epigraphs). For a function  $f: V \to (-\infty, \infty]$ , its *epigraph* is defined as

$$\operatorname{epi}(f) := \{ (\boldsymbol{x}, t) \in \boldsymbol{V} \times \mathbb{R} : t \ge f(\boldsymbol{x}) \}$$

Note that

$$\left\{ (\boldsymbol{x}, f(\boldsymbol{x})) : \boldsymbol{x} \in \operatorname{dom}(f) \right\} \subset \operatorname{epi}(f) \subset \operatorname{dom}(f) \times \mathbb{R}$$

A useful fact is that  $f : \mathbf{V} \to (-\infty, \infty]$  is convex if and only if epi(f) is a convex subset of  $\mathbf{V} \times \mathbb{R}$ .

Indeed, suppose first that f is a convex function. Let  $(x, s), (y, t) \in epi(f)$ . Then for  $\lambda \in (0, 1)$ ,

$$(1-\lambda)s + \lambda t \geq (1-\lambda)f(\boldsymbol{x}) + \lambda f(\boldsymbol{y}) \geq f((1-\lambda)\boldsymbol{x} + \lambda \boldsymbol{y})$$

and therefore

$$(1-\lambda)(\boldsymbol{x},s) + \lambda(\boldsymbol{y},t) = ((1-\lambda)\boldsymbol{x} + \lambda\boldsymbol{y}, (1-\lambda)s + \lambda t) \in \operatorname{epi}(f).$$

This shows convexity of epi(f).

Assume now that epi(f) is a convex set. For  $\boldsymbol{x}, \boldsymbol{y} \in dom(f)$  and  $\lambda \in (0, 1)$ ,

$$\left((1-\lambda)\boldsymbol{x}+\lambda\boldsymbol{y},\,(1-\lambda)f(\boldsymbol{x})+\lambda f(\boldsymbol{y})\right) = (1-\lambda)(\boldsymbol{x},f(\boldsymbol{x}))+\lambda(\boldsymbol{y},f(\boldsymbol{y})) \in \operatorname{epi}(f),$$

because  $(\boldsymbol{x}, f(\boldsymbol{x})), (\boldsymbol{y}, f(\boldsymbol{y})) \in \operatorname{epi}(f)$  and  $\operatorname{epi}(f)$  is convex. This shows that

$$(1-\lambda)f(\boldsymbol{x}) + \lambda f(\boldsymbol{y}) \geq f((1-\lambda)\boldsymbol{x} + \lambda \boldsymbol{y}).$$

Hence f is convex.

Convexity of a function can often be established by means of the following lemma the proof of which is left to the reader as an exercise:

Lemma 3.20 (Verifying convexity).

(a) Let  $m \in \mathbb{N}$ , and let  $f := \sum_{i=1}^{m} a_i f_i$  with numbers  $a_i > 0$  and convex functions  $f_i : \mathbf{V} \to (-\infty, \infty]$ . Show that f is convex, too.

(b) Let L be an affine function from V into another real vector space W, and let  $g : W \to (-\infty, \infty]$  be a convex function. Then  $g \circ L$  is a convex function on V.

(c) Let  $(f_{\lambda})_{\lambda \in \Lambda}$  be a family of convex functions  $f_{\lambda} : V \to (-\infty, \infty]$ . Then the pointwise supremum  $\sup_{\lambda \in \Lambda} f_{\lambda}$  is convex, too.

Here is a simple but very useful observation: A function  $f : C \to (-\infty, \infty]$  is convex if and only if for arbitrary points  $x \in C$  and vectors  $v \in V \setminus \{0\}$  the function

$$t \mapsto f(\boldsymbol{x} + t\boldsymbol{v})$$

is convex on the interval  $\{t \in \mathbb{R} : x + tv \in C\}$ . Analogously, a function  $f : C \to \mathbb{R}$  is strictly concave, if the previous univariate function is always strictly convex.

**Exercise 3.21** (Affine functions). Let  $f : \mathbf{V} \to \mathbb{R}$  be both convex and concave. Show that f is an *affine* function, that means,

$$f(\boldsymbol{x}) = c + L(\boldsymbol{x})$$

for some constant  $c \in \mathbb{R}$  and a linear function  $L : \mathbf{V} \to \mathbb{R}$ .

**Lemma 3.22** (Monotone gradients). Let C be an open, convex subset of  $\mathbb{R}^d$ , and let  $f : C \to \mathbb{R}$  be differentiable. Then f is convex (strictly convex) if and only if

$$\left( 
abla f(oldsymbol{y}) - 
abla f(oldsymbol{x}) 
ight)^{ op} (oldsymbol{y} - oldsymbol{x}) \ge (>) \ 0$$

for arbitrary different points  $x, y \in C$ .

60

**Lemma 3.23** (Hessian matrices). Let C be an open, convex subset of  $\mathbb{R}^d$ , and let  $f : C \to \mathbb{R}$  be twice differentiable.

(a) f is convex if and only if the Hessian matrix  $D^2 f(x)$  is positive semidefinite for any point  $x \in C$ .

(b) f is strictly convex if the Hessian matrix  $D^2 f(x)$  is positive definite for any point  $x \in C$ .

The previous two lemmata follow essentially from our results about univariate functions (Section 3.1) and the observation that for any  $x \in C$  and  $v \in \mathbb{R}^d$ , g(t) := f(x + tv) satisfies

$$g'(t) = \nabla f(\boldsymbol{x} + t\boldsymbol{v})^{\top} \boldsymbol{v}$$
 and  $g''(t) = \boldsymbol{v}^{\top} D^2 f(\boldsymbol{x} + t\boldsymbol{v}) \boldsymbol{v}.$ 

**Example 3.24** (Quadratic functions). Let  $f : \mathbb{R}^d \to \mathbb{R}$  be a quadratic function, i.e.  $f(\boldsymbol{x}) = c + \boldsymbol{b}^\top \boldsymbol{x} + 2^{-1} \boldsymbol{x}^\top \boldsymbol{A} \boldsymbol{x}$  with a constant  $c \in \mathbb{R}$ , a vector  $\boldsymbol{b} \in \mathbb{R}^d$  and a symmetric matrix  $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ . Here

$$\nabla f(\boldsymbol{x}) = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}$$
 and  $D^2 f(\boldsymbol{x}) = \boldsymbol{A}$ .

Note also that

$$\langle \nabla f(\boldsymbol{y}) - \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle = (\boldsymbol{y} - \boldsymbol{x})^{\top} \boldsymbol{A}(\boldsymbol{y} - \boldsymbol{x})$$

Hence f is convex if, and only, if A is positiv semidefinite. Strict convexity of f is equivalent to A being positive definite.

A simple version of Jensen's inequality. Inductively one can show that a convex function f on C has the following property: For any  $n \in \mathbb{N}$ ,

$$f\Big(\sum_{i=1}^n \lambda_i oldsymbol{x}_i\Big) \ \le \ \sum_{i=1}^n \lambda_i f(oldsymbol{x}_i)$$

whenever  $x_1, \ldots, x_n \in C$  and  $\lambda_1, \ldots, \lambda_n > 0$  such that  $\sum_{i=1}^n \lambda_i = 1$ . This inequality is sometimes called Jensen's inequality. It implies a simple fact which we shall use several times:

**Lemma 3.25.** Let  $C = \operatorname{conv}(S)$  for a nonvoid set  $S \subset V$ . For any convex function  $f : C \to (-\infty, \infty]$ ,

$$\sup_{\boldsymbol{x}\in\boldsymbol{C}} f(\boldsymbol{x}) = \sup_{\boldsymbol{x}\in\boldsymbol{S}} f(\boldsymbol{x}).$$

For any affine function  $f : \mathbf{C} \to \mathbb{R}$ ,

$$\inf_{\boldsymbol{x}\in \boldsymbol{C}} f(\boldsymbol{x}) = \inf_{\boldsymbol{x}\in \boldsymbol{S}} f(\boldsymbol{x}) \text{ and } \sup_{\boldsymbol{x}\in \boldsymbol{C}} f(\boldsymbol{x}) = \sup_{\boldsymbol{x}\in \boldsymbol{S}} f(\boldsymbol{x}).$$

**Proof of Lemma 3.25.** Since  $C = \operatorname{conv}(S)$ ,  $\sup_{x \in C} f(x)$  is the supremum of

$$f\left(\sum_{i=1}^n \lambda_i \boldsymbol{x}_i\right)$$

over all  $n \in \mathbb{N}$ ,  $x_1, \ldots, x_n \in S$  and  $\lambda_1, \ldots, \lambda_n > 0$  with  $\sum_{i=1}^n \lambda_i = 1$ . Taking n = 1 shows that  $\sup_{x \in C} f(x) \ge \sup_{x \in S} f(x)$ . But by Jensen's inequality,

$$f\Big(\sum_{i=1}^n \lambda_i \boldsymbol{x}_i\Big) \ \le \ \sum_{i=1}^n \lambda_i f(\boldsymbol{x}_i) \ \le \ \sup_{\boldsymbol{x} \in \boldsymbol{S}} f(\boldsymbol{x}).$$

If f is an affine function, then f as well as -f are convex, so the second part is an immediate consequence of the first part.

**Exercise 3.26** (Minkowski inequality). Prove the generalized Minkowski inequality: For arbitrary positive numbers  $x_1, \ldots, x_n$  and  $\lambda_1, \ldots, \lambda_n$  with  $\sum_{i=1}^n \lambda_i = 1$ ,

$$\prod_{i=1}^{n} x_i \leq \sum_{i=1}^{n} \lambda_i x_i^{1/\lambda_i}$$

**Exercise 3.27** (Investment plans). Mr. Müller and Mr. Lüdenscheid, two wealthy gentlemen, invest money at time points  $t_1 < t_2 < \cdots < t_n$  in shares of the promising company *IMSV*. The strategy of Mr. Müller is to invest each time a fixed amount c > 0, regardless of the current price. (For simplicity we assume that one can buy fractions of shares.) Mr. Lüdenscheid, however, buys each time a fixed number d > 0 of shares. Who of the two will end up with the higher ratio

amount of shares amount of invested money?

## **3.3** Directional Derivatives and Minimizers

One of the most important features of convex functions is the fact that any local minimizer is also a global minimizer. This will be derived subsequently.

Let  $f : V \to (-\infty, \infty]$  be a convex function. For any point  $x \in \text{dom}(f)$  and arbitrary vectors  $v \in V$ , the one-sided directional derivative

$$Df(\boldsymbol{x}, \boldsymbol{v}) := \lim_{t \to 0+} rac{f(\boldsymbol{x} + t \boldsymbol{v}) - f(\boldsymbol{x})}{t} \in [-\infty, \infty]$$

is well-defined. This follows from Corollary 3.6, applied to the convex function  $\mathbb{R} \ni t \mapsto f(\boldsymbol{x} + t\boldsymbol{v})$ . As a function of  $\boldsymbol{v}, Df(\boldsymbol{x}, \boldsymbol{v})$  is sublinear, that means,

- (3.5)  $Df(\boldsymbol{x},\lambda\boldsymbol{v}) = \lambda Df(\boldsymbol{x},\boldsymbol{v})$  for  $\boldsymbol{v} \in \boldsymbol{V}$  and  $\lambda \ge 0$ ,
- $(3.6) \quad Df(\boldsymbol{x}, \boldsymbol{v} + \boldsymbol{w}) \leq Df(\boldsymbol{x}, \boldsymbol{v}) + Df(\boldsymbol{x}, \boldsymbol{w}) \quad \text{ for } \boldsymbol{v}, \boldsymbol{w} \in \boldsymbol{V} \text{ such that} \\ \left\{ Df(\boldsymbol{x}, \boldsymbol{v}), Df(\boldsymbol{x}, \boldsymbol{w}) \right\} \neq \{-\infty, \infty\}.$

In particular, if  $Df(x, v) > -\infty$  for all  $v \in V$ , then  $Df(x, \cdot)$  is a convex function on V.

**Exercise 3.28**. Prove the statements (3.5) and (3.6). Show also that  $Df(x, \cdot)$  is a convex function on V, provided that  $Df(x, \cdot) > -\infty$ .

**Example 3.29**. To see that for fixed  $x \in dom(f)$ ,  $Df(x, \cdot)$  may take the values  $\pm \infty$ , let  $V = \mathbb{R}$  and

$$f(x) := \begin{cases} \infty & \text{if } |x| > 1, \\ -\sqrt{1 - x^2} & \text{if } |x| \le 1. \end{cases}$$

Then  $Df(x,v) = vf'(x) = vx/\sqrt{1-x^2}$  for  $x \in (-1,1)$  and  $v \in \mathbb{R}$ , but  $Df(1,\pm 1) = \pm \infty$  and  $Df(-1,\pm 1) = \mp \infty$ .

It follows from Corollary 3.6 that

(3.7) 
$$f(\boldsymbol{x}) + Df(\boldsymbol{x}, \boldsymbol{v}) \leq f(\boldsymbol{x} + \boldsymbol{v}) \quad \text{for all } \boldsymbol{v} \in \boldsymbol{V}.$$

This statement and the definition of directional derivatives imply the aforementioned fact that every local minimizer of a convex function is automatically a global minimizer:

**Theorem 3.30.** Let  $f : \mathbf{V} \to (-\infty, \infty]$  be a convex function. The following two statements about  $\mathbf{x} \in \text{dom}(f)$  are equivalent:

- (a) x is a global minimizer of f, that means,  $f(x) \le f(y)$  for all  $y \in V$ ;
- (b)  $Df(\boldsymbol{x}, \boldsymbol{v}) \geq 0$  for arbitrary  $\boldsymbol{v} \in \boldsymbol{V}$ .

Here is a version of the same result for convex functions on convex subsets of V:

**Theorem 3.31**. Let *f* be a real-valued function on a convex subset *C* of *V*. The following two statements about  $x \in C$  are equivalent:

(a) x is a global minimizer of f on C, that means,  $f(x) \leq f(y)$  for all  $y \in C$ ;

(b)  $Df(\boldsymbol{x}, \boldsymbol{y} - \boldsymbol{x}) \geq 0$  for arbitrary  $\boldsymbol{y} \in \boldsymbol{C}$ .

**Exercise 3.32** ( $\ell_p$ -norms). For  $p \ge 1$  and  $\boldsymbol{x} \in \mathbb{R}^d$  let  $\|\boldsymbol{x}\|_p := \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$ . Determine for  $p, q \ge 1$  the number

$$C_{p,q} := \max_{\boldsymbol{x} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}} rac{\|\boldsymbol{x}\|_q}{\|\boldsymbol{x}\|_p}.$$

Hint: One can write  $(||\boldsymbol{x}||_q/||\boldsymbol{x}||_p)^q$  as a symmetric convex or concave function of the vector  $\boldsymbol{w} := ||\boldsymbol{x}||_p^{-p}(|\boldsymbol{x}_i|^p)_{i=1}^d$  which lies in the convex polyhedron  $\boldsymbol{C} = \{\boldsymbol{w} \in [0,\infty)^d : \sum_{i=1}^d w_i = 1\}$ . **Example 3.33** (Differentiable and quadratic functions). Let  $\boldsymbol{C}$  be an open subset of  $\mathbb{R}^d$ , and let  $f : \boldsymbol{C} \to \mathbb{R}$  be differentiable and convex. Since  $Df(\boldsymbol{x}, \boldsymbol{v}) = \nabla f(\boldsymbol{x})^\top \boldsymbol{v}$  for arbitrary  $\boldsymbol{x} \in \boldsymbol{C}$  and  $\boldsymbol{v} \in \mathbb{R}^d$ , a point  $\boldsymbol{x} \in \boldsymbol{C}$  minimizes f if and only if

$$\nabla f(\boldsymbol{x}) = \boldsymbol{0}.$$

In particular, let  $C = \mathbb{R}^d$ , and let f be a quadratic function, i.e.  $f(x) = c + b^\top x + 2^{-1} x^\top A x$ for a constant  $c \in \mathbb{R}$ , a vector  $b \in \mathbb{R}^d$  and a symmetric, positive semidefinite matrix  $A \in \mathbb{R}^{d \times d}$ . Since  $\nabla f(x) = Ax + b$ , a point x minimizes f if and only if

$$Ax = -b.$$

In particular, a minimizer exists if and only if b lies in the column space of A. In case of A being positive definite, the function f is strictly convex with unique minimizer

$$\boldsymbol{x} = -\boldsymbol{A}^{-1}\boldsymbol{b}.$$

**Exercise 3.34**. For  $\boldsymbol{x} \in (0, \infty)^2$  let

$$f(\boldsymbol{x}) := \frac{x_1^2}{x_2} + \frac{x_2^2}{x_1} - x_1 - x_2.$$

Show that  $f \ge 0$ , and determine the set of minimizers of f.

**Remark 3.35**. Let f be a convex function on an open convex subset C of  $\mathbb{R}^d$ , and let  $e_1, \ldots, e_d$  be the standard basis of  $\mathbb{R}^d$ . In case of  $d \ge 2$ , it may happen that a point  $x \in C$  satisfies the 2d constraints  $Df(x, \pm e_i) \ge 0, 1 \le i \le d$ , but x is not a minimizer of f. As a counterexample, consider

$$f(\boldsymbol{x}) := \max_{1 \leq i \leq d} x_i.$$

Since f is a maximum of d linear functions, it is convex. Moreover, one can show that

$$Df(\boldsymbol{x}, \boldsymbol{v}) = \max_{i \in I(\boldsymbol{x})} v_i$$
 with  $I(\boldsymbol{x}) := \operatorname*{arg\,max}_{i \in \{1, ..., d\}} x_i$ 

In particular, if  $x = r\mathbf{1}_d$ , then for  $1 \le i \le d$ ,

$$Df(\boldsymbol{x}, \boldsymbol{e}_i) = 1$$
 and  $Df(\boldsymbol{x}, -\boldsymbol{e}_i) = 0$ .

But  $Df(\boldsymbol{y}, -\mathbf{1}_d) = -1$  for any  $\boldsymbol{y} \in \mathbb{R}^d$ , whence f has no minimizer on any open set  $\boldsymbol{C} \subset \mathbb{R}^d$ .

**Exercise 3.36.** Find a convex function  $f : \mathbb{R}^d \to \mathbb{R}$  such that for a given basis  $b_1, \ldots, b_d$  of  $\mathbb{R}^d$  and some  $x \in \mathbb{R}^d$ , the 2*d* inequalites  $Df(x, \pm b_i) > 0, 1 \le i \le d$ , are satisfied but  $x \notin \arg \min_{\mathbb{R}^d} (f)$ .

**Exercise 3.37** (Spatial median, I). For given points  $x_1, \ldots, x_n \in \mathbb{R}^d$  we consider the function  $f : \mathbb{R}^d \to \mathbb{R}$  with

$$f(x) := \sum_{i=1}^{n} \|x_i - x\|.$$

(a) Show that f is a convex function satisfying  $f(x) \to \infty$  as  $||x|| \to \infty$ .

(b) Derive an explicit expression for the directional derivatives Df(x, v). (You should distinguish the index sets  $\{i : x_i = x\}$  and  $\{i : x_i \neq x\}$ .)

(c) Show that

$$rgmin_{oldsymbol{x}\in\mathbb{R}^d} f(oldsymbol{x}) \ \subset \ \mathrm{conv}(oldsymbol{x}_1,\ldots,oldsymbol{x}_n).$$

#### Two examples of regularized estimators

**Example 3.38** (LASSO penalization). In statistics one is often interested in minimizing a differentiable convex function  $g : \mathbb{R}^d \to \mathbb{R}$ , for instance a sum of squared residuals

$$g(\boldsymbol{x}) = \sum_{i=1}^{n} (y_i - \boldsymbol{z}_i^{\top} \boldsymbol{x})^2 / 2$$

with a given data vector  $\boldsymbol{y} \in \mathbb{R}^n$  and given design vectors  $\boldsymbol{z}_1, \ldots, \boldsymbol{z}_n \in \mathbb{R}^d$ . In this particular example for g, if the vectors  $\boldsymbol{z}_i$  span  $\mathbb{R}^d$ , there exists a unique minimizer  $\hat{\boldsymbol{x}}$  of g: With  $\boldsymbol{Z} := [\boldsymbol{z}_1, \ldots, \boldsymbol{z}_n]^\top \in \mathbb{R}^{n \times d}$  and  $\boldsymbol{\Gamma} := \boldsymbol{Z}^\top \boldsymbol{Z} = \sum_{i=1}^n \boldsymbol{z}_i \boldsymbol{z}_i^\top$ ,

$$g(\boldsymbol{x}) = \|\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{x}\|^2/2$$
  
=  $\|\boldsymbol{y}\|^2/2 - \boldsymbol{x}^\top \boldsymbol{Z}^\top \boldsymbol{y} + \boldsymbol{x}^\top \boldsymbol{\Gamma}\boldsymbol{x}/2$   
=  $\|\boldsymbol{y}\|^2/2 - \boldsymbol{y}^\top \boldsymbol{Z} \boldsymbol{\Gamma}^{-1} \boldsymbol{Z}^\top \boldsymbol{y}/2 + (\boldsymbol{x} - \boldsymbol{\Gamma}^{-1} \boldsymbol{Z}^\top \boldsymbol{y})^\top \boldsymbol{\Gamma} (\boldsymbol{x} - \boldsymbol{\Gamma}^{-1} \boldsymbol{Z}^\top \boldsymbol{y})/2$   
 $\geq \|\boldsymbol{y}\|^2/2 - \boldsymbol{y}^\top \boldsymbol{Z} \boldsymbol{\Gamma}^{-1} \boldsymbol{Z}^\top \boldsymbol{y}/2$ 

with equality if and only if

$$\boldsymbol{x} = \boldsymbol{\Gamma}^{-1} \boldsymbol{Z}^{ op} \boldsymbol{y}.$$

But often one is interested in computing parameter vectors x having only a few nonzero components while g(x) is still relatively small. In the least squares example it may even happen that d > n, so the minimization problem has no unique solution. One way to deal with these problems is to minimize the penalized function

$$f_{\lambda}(\boldsymbol{x}) := g(\boldsymbol{x}) + \lambda \sum_{j=1}^{d} |x_j|.$$

for some tuning parameter  $\lambda > 0$ . (This method is called least absolute shrinkage and selection operator.) Here one can easily show that

$$Df_{\lambda}(\boldsymbol{x}, \boldsymbol{v}) = \sum_{j=1}^{d} \left( \frac{\partial g}{\partial x_{j}}(\boldsymbol{x})v_{j} + \lambda \operatorname{sign}(x_{j})v_{j} + \lambda 1_{[x_{j}=0]}|v_{j}| \right)$$
$$= \sum_{j:x_{j}\neq 0} \left( \frac{\partial g}{\partial x_{j}}(\boldsymbol{x}) + \lambda \operatorname{sign}(x_{j}) \right)v_{j} + \sum_{j:x_{j}=0} \left( \frac{\partial g}{\partial x_{j}}(\boldsymbol{x})v_{j} + \lambda |v_{j}| \right).$$

From this representation one can easily deduce that a vector  $\widehat{x} \in \mathbb{R}^d$  minimizes  $f_{\lambda}$  if and only if

$$\frac{\partial g}{\partial x_j}(\widehat{\boldsymbol{x}}) \begin{cases} \in [-\lambda,\lambda] & \text{if } \widehat{x}_j = 0, \\ = -\lambda \operatorname{sign}(\widehat{x}_j) & \text{if } \widehat{x}_j \neq 0. \end{cases}$$

In particular,  $\hat{x}_j = 0$  whenever the modulus of  $\partial g(\hat{x}) / \partial x_j$  is less than  $\lambda$ . This explains why the LASSO approach tends to produce vectors  $\hat{x}$  with several components equal to zero.

**Exercise 3.39.** Let  $g : \mathbb{R}^d \to \mathbb{R}$  be given by  $g(\boldsymbol{x}) = \|\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{x}\|^2/2$  with given numbers  $y_1, \ldots, y_n \in \mathbb{R}$  and a matrix  $\boldsymbol{Z} \in \mathbb{R}^{n \times d}$  such that  $\sum_{i=1}^n Z_{ij}^2 = 1$  for  $1 \le j \le d$ . For  $\lambda > 0$  we consider  $f_{\lambda}(\boldsymbol{x}) := g(\boldsymbol{x}) + \lambda \sum_{j=1}^d |x_j|$ .

(a) Determine for  $x \in \mathbb{R}^d$  and any  $k \in \{1, \ldots, d\}$  the (unique!) vector

$$\boldsymbol{\psi}^{(k)}(\boldsymbol{x}) := \arg\min\Big\{f_{\lambda}(\tilde{\boldsymbol{x}}): \tilde{\boldsymbol{x}} \in \mathbb{R}^d, \tilde{x}_j = x_j \text{ whenever } j \neq k\Big\}.$$

Show that this defines a continuous mapping  $\psi^{(k)} : \mathbb{R}^d \to \mathbb{R}^d$ .

These mappings  $\psi^{(k)}$  are the basis of a simple iterative algorithm to minimize  $f_{\lambda}$ : For a given candidate x we compute  $x_{\text{new}}$  as follows: After initializing  $x_{\text{new}} \leftarrow x$ , we execute  $x_{\text{new}} \leftarrow \psi^{(k)}(x_{\text{new}})$  for k = 1, ..., d. As long as the difference  $f_{\lambda}(x) - f_{\lambda}(x_{\text{new}})$  is larger than a given small constant  $\delta > 0$ , we set  $x \leftarrow x_{\text{new}}$  and repeat this loop. A justification of this algorithm will be given in the next chapter.

(b) Suppose that  $x = \psi^{(k)}(x)$  for  $1 \le k \le d$ . Show that x is a minimizer of  $f_{\lambda}$ .

**Example 3.40** (Taut Strings). Consider a vector  $\boldsymbol{y} \in \mathbb{R}^n$  of observations  $y_1, y_2, \ldots, y_n$  corresponding, for instance, to measurements at *n* consecutive time points. If these measurements are noisy, one would like to approximate  $\boldsymbol{y}$  by a 'smoother' vector  $\boldsymbol{x}$ .

A classical approach is to minimize for a given tuning parameter  $\lambda > 0$  the target function

$$f_{\lambda}(\boldsymbol{x}) := \sum_{i=1}^{n} (y_i - x_i)^2 + \lambda \sum_{j=1}^{n-1} (x_{j+1} - x_j)^2$$
$$= \|\boldsymbol{y} - \boldsymbol{x}\|^2 + \lambda \boldsymbol{x}^{\top} \boldsymbol{S} \boldsymbol{x}$$

with the symmetric and positive semidefinite matrix  $S \in \mathbb{R}^{n \times n}$  given by

$$S_{ij} := \begin{cases} 1 & \text{if } i = j \in \{1, n\}, \\ 2 & \text{if } 1 < i = j < n, \\ 1 & \text{if } |i - j| = 1, \\ 0 & \text{else.} \end{cases}$$

The first summand  $\|\boldsymbol{y} - \boldsymbol{x}\|^2$  of  $f_{\lambda}(\boldsymbol{x})$  measures the distance between  $\boldsymbol{x}$  and  $\boldsymbol{y}$  while the second summand  $\lambda \boldsymbol{x}^{\top} \boldsymbol{S} \boldsymbol{x}$  quantifies how strongly consecutive components of  $\boldsymbol{x}$  differ. This function  $f_{\lambda}$  has a unique minimizer  $\hat{\boldsymbol{x}}_{\lambda}$ :

$$egin{aligned} f_\lambda(oldsymbol{x}) &= \|oldsymbol{y}\|^2 - 2oldsymbol{x}^ opoldsymbol{y} + oldsymbol{x}^ op(oldsymbol{I}_n + \lambdaoldsymbol{S})^{-1}oldsymbol{y} \ &= \|oldsymbol{y}\|^2 - oldsymbol{y}^ op(oldsymbol{I}_n + \lambdaoldsymbol{S})^{-1}oldsymbol{y} \ &+ oldsymbol{x} - (oldsymbol{I}_n + \lambdaoldsymbol{S})^{-1}oldsymbol{y} \ &\geq \|oldsymbol{y}\|^2 - oldsymbol{y}^ op(oldsymbol{I}_n + \lambdaoldsymbol{S})^{-1}oldsymbol{y} \ &\leq \|oldsymbol{S}\| oldsymbol{S}^ op(oldsymbol{S})^ op(ol$$

with equality if and only if

$$\boldsymbol{x} = (\boldsymbol{I}_n + \lambda \boldsymbol{S})^{-1} \boldsymbol{y}.$$

Let us now consider an alternative approach: We want to minimize

$$f_{\lambda}(\boldsymbol{x}) := \|\boldsymbol{y} - \boldsymbol{x}\|^2 / 2 + \lambda \sum_{j=1}^{n-1} |x_{j+1} - x_j|$$

The latter penalty term favours vectors  $\boldsymbol{x}$  such that  $x_j = x_{j+1}$  for several indices j < n. Note first that  $f_{\lambda}$  is a strictly convex function such that  $f_{\lambda}(\boldsymbol{x}) \to \infty$  as  $\|\boldsymbol{x}\| \to \infty$ . This implies that  $f_{\lambda}$  has a unique minimizer  $\hat{\boldsymbol{x}}_{\lambda}$ , see also Section 3.5.3. Now we characterize  $\hat{\boldsymbol{x}}_{\lambda}$  by means of one-sided directional derivatives. One can easily show that for arbitrary vectors  $\boldsymbol{x}, \boldsymbol{v} \in \mathbb{R}^n$ ,

$$Df_{\lambda}(\boldsymbol{x}, \boldsymbol{v}) = (\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{v} + \lambda \sum_{j=1}^{n-1} (\operatorname{sign}(x_{j+1} - x_j)(v_{j+1} - v_j) + 1_{[x_j = x_{j+1}]} |v_{j+1} - v_j|).$$

This does not look too illuminating at first glance, but plugging in some special vectors v reveals interesting information about  $\hat{x}_{\lambda}$ . Specifically, let  $b_k := (1_{[i \le k]})_{i=1}^n$  for k = 1, 2, ..., n.

If 
$$v = \pm b_n$$
, then

$$Df_{\lambda}(\boldsymbol{x}, \boldsymbol{v}) = \pm (\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{b}_{n} = \pm \Big(\sum_{i=1}^{n} x_{i} - \sum_{i=1}^{n} y_{i}\Big).$$

Hence a first necessary condition for  $oldsymbol{x} = \widehat{oldsymbol{x}}_{\lambda}$  is

(3.8) 
$$\sum_{i=1}^{n} x_i = \sum_{i=1}^{n} y_i$$

Next consider  $v = \pm b_k$  for some index  $1 \le k < n$ . Then

$$Df_{\lambda}(\boldsymbol{x}, \boldsymbol{v}) = \pm (\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{b}_{k} \mp \lambda \operatorname{sign}(x_{k+1} - x_{k}) + \lambda \mathbf{1}_{[x_{k} = x_{k+1}]}$$
$$= \pm \left(\sum_{i=1}^{k} x_{i} - \sum_{i=1}^{k} y_{i}\right) \mp \lambda \operatorname{sign}(x_{k+1} - x_{k}) + \lambda \mathbf{1}_{[x_{k} = x_{k+1}]}.$$

This leads to further necessary conditions for  $\boldsymbol{x} = \widehat{\boldsymbol{x}}_{\lambda}$ : For  $1 \leq k < n$ ,

(3.9) 
$$\left|\sum_{i=1}^{k} x_i - \sum_{i=1}^{k} y_i\right| \leq \lambda,$$

(3.10) 
$$\sum_{i=1}^{k} x_i = \sum_{i=1}^{k} y_i + \lambda \operatorname{sign}(x_{k+1} - x_k) \quad \text{if } x_k \neq x_{k+1}.$$

In fact, Conditions (3.8), (3.9) and (3.10) are even sufficient for  $x = \hat{x}_{\lambda}$ . This can be verified by writing an arbitrary vector v as

$$oldsymbol{v} \;=\; \sum_{k=1}^{n-1} (v_k - v_{k+1}) oldsymbol{b}_k + v_n oldsymbol{b}_n$$

and noting that

$$Df_{\lambda}(\boldsymbol{x}, \boldsymbol{v}) = \sum_{k=1}^{n-1} (v_{k} - v_{k+1})(\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{b}_{k} + v_{n} \underbrace{(\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{b}_{n}}_{=0} \\ + \lambda \sum_{j=1}^{n-1} (\operatorname{sign}(x_{j+1} - x_{j})(v_{j+1} - v_{j}) + 1_{[x_{j} = x_{j+1}]} |v_{j+1} - v_{j}|) \\ = \sum_{k=1}^{n-1} (v_{k} - v_{k+1})(\boldsymbol{x} - \boldsymbol{y})^{\top} \boldsymbol{b}_{k} \\ + \lambda \sum_{k=1}^{n-1} (-(v_{k} - v_{k+1}) \operatorname{sign}(x_{k+1} - x_{k}) + 1_{[x_{j} = x_{j+1}]} |v_{k} - v_{k+1}|) \\ = \sum_{k=1}^{n-1} |v_{k} - v_{k+1}| Df_{\lambda}(\boldsymbol{x}, \operatorname{sign}(v_{k} - v_{k+1}) \boldsymbol{b}_{k})$$

is nonnegative.

There is a nice geometrical interpretation of Conditions (3.8), (3.9) and (3.10): Any vector  $w \in \mathbb{R}^n$  can be identified with its partial sum function,

$$\{0, 1, \dots, n\} \ni k \mapsto S_{\boldsymbol{w}}(k) := \begin{cases} 0 & \text{if } k = 0, \\ \sum_{i=1}^{k} w_i & \text{if } k > 0. \end{cases}$$

Condition (3.8) is equivalent to saying that  $S_x$  and  $S_y$  coincide at 0 and n. Condition (3.9) means that  $S_y - \lambda \leq S_x \leq S_y + \lambda$ . Condition (3.10) means the following: If  $S_x$  bends upwards at an index  $k \in \{1, \ldots, n-1\}$  (i.e.  $x_k < x_{k+1}$ ), then it hits the upper boundary  $S_y + \lambda$  there; if it bends downwards, then it hits the lower boundary  $S_y - \lambda$ .

This is illustrated in Figures 3.1, 3.2 and 3.3. For n = 70 and  $\lambda = 2, 1, 0.5$ , respectively, the upper panel shows a scatter plot of the pairs  $(i, y_i)$  (black bullets) and  $(i, \hat{x}_{\lambda,i})$  (blue circles and line). In the lower panel one sees the partial sum function  $S_y$  (thin black line) and the same function  $\pm \lambda$  on  $\{1, \ldots, n-1\}$  (black lines and bullets). In between one sees the partial sum function  $S_{\hat{x}_{\lambda}}$  (thick blue line). All these functions are interpolated linearly between consecutive integers. Obviously the fitted vector  $\hat{x}_{\lambda}$  gets closer to y if  $\lambda$  is lowered.

When looking at the lower panels of these figures, one can imagine the following procedure: Suppose the coordinate system is a wooden board, and each bullet represents a nail. Now we connect the leftmost and rightmost nail with a string lying loosely between the lower and upper row of nails. Now we pull the string tight and obtain the graph of the partial sum function of  $\hat{x}_{\lambda}$ . This explains the name 'taut string' method. The actual computation is possible in O(n) steps, but the precise algorithm is beyond the scope of this lecture.

## **3.4 Smoothness Properties**

In this section we show that convex functions have almost automatically certain smoothness properties. Throughout this section let  $(V, \|\cdot\|)$  be a real normed vector space with closed unit ball B, and let  $f : V \to (-\infty, \infty]$  be a convex function.

### 3.4.1 Local Lipschitz continuity

The next two results are about local Lipschitz continuity of convex functions.

**Theorem 3.41** (Local Lipschitz continuity, I). Suppose that for some  $x \in V$  and  $\epsilon > 0$  the function f is bounded from above on the set  $x + \epsilon B$ . Then f is Lipschitz continuous on  $x + \delta B$  whenever  $0 < \delta < \epsilon$ .

**Theorem 3.42** (Local Lipschitz continuity, II). Let  $V = \mathbb{R}^d$ . Then f is continuous on the interior of dom(f). It is even Lipschitz continuous on any compact subset of the interior of dom(f).

**Proof of Theorem 3.41.** For two different points  $\boldsymbol{y}, \boldsymbol{z} \in \text{dom}(f)$  we consider the mapping  $\tilde{f}$ :  $\mathbb{R} \to (-\infty, \infty]$  with  $\tilde{f}(t) := f(\boldsymbol{y} + t(\boldsymbol{z} - \boldsymbol{y}))$  and the unit vector  $\boldsymbol{v} := \|\boldsymbol{z} - \boldsymbol{y}\|^{-1}(\boldsymbol{z} - \boldsymbol{y})$ . Then

$$f(\boldsymbol{z}) - f(\boldsymbol{y}) \begin{cases} \leq \tilde{f}'(1+) = + \|\boldsymbol{z} - \boldsymbol{y}\| Df(\boldsymbol{z}, + \boldsymbol{v}), \\ \geq \tilde{f}'(0-) = - \|\boldsymbol{z} - \boldsymbol{y}\| Df(\boldsymbol{y}, - \boldsymbol{v}). \end{cases}$$

Hence it suffices to show that

$$\sup_{\boldsymbol{y} \in \boldsymbol{x} + \delta \boldsymbol{B}, \boldsymbol{v} \in \boldsymbol{B}} Df(\boldsymbol{y}, \boldsymbol{v}) < \infty$$

for  $0 < \delta < \epsilon$ .

To this end let  $M := \sup_{x+\epsilon B} (f) < \infty$ . For any  $y \in x + \epsilon B$ , the vector  $\tilde{y} = x - (y - x)$ belongs to  $x + \epsilon B$ , too, so it follows from  $x = (y + \tilde{y})/2$  that  $(f(y) + f(\tilde{y}))/2 \ge f(x)$ , whence

$$f(\boldsymbol{y}) \geq 2f(\boldsymbol{x}) - f(\widetilde{\boldsymbol{y}}) \geq 2f(\boldsymbol{x}) - M$$



Figure 3.1: Illustration of the Taut String method,  $\lambda = 2$ .



Figure 3.2: Illustration of the Taut String method,  $\lambda = 1$ .



Figure 3.3: Illustration of the Taut String method,  $\lambda=0.5.$ 

This shows that

$$\inf_{\boldsymbol{x}+\boldsymbol{\epsilon}\boldsymbol{B}} f \geq 2f(\boldsymbol{x}) - M.$$

Hence for  $\boldsymbol{y} \in \boldsymbol{x} + \delta \boldsymbol{B}$  and  $\boldsymbol{v} \in \boldsymbol{B}$ ,

$$Df(\boldsymbol{y}, \boldsymbol{v}) \leq \frac{f(\boldsymbol{y} + (\epsilon - \delta)\boldsymbol{v}) - f(\boldsymbol{y})}{\epsilon - \delta} \leq \frac{M - (2f(\boldsymbol{x}) - M)}{\epsilon - \delta} = \frac{2(M - f(\boldsymbol{x}))}{\epsilon - \delta}.$$

**Proof of Theorem 3.42.** Since all norms on  $\mathbb{R}^d$  are equivalent, we consider temporarily the norm  $||\boldsymbol{v}|| := \max_i |v_i|$  with the corresponding unit ball  $\boldsymbol{B} = [-1, 1]^d$ . Now let  $\boldsymbol{x}$  be an interior point of dom(f). Then there exists an  $\epsilon = \epsilon(\boldsymbol{x}) > 0$  such that  $\boldsymbol{x} + \epsilon \boldsymbol{B} \subset \text{dom}(f)$ . Since the cube  $\boldsymbol{x} + \epsilon \boldsymbol{B}$  is the convex hull of the finite set

$$\{\boldsymbol{x} + \boldsymbol{\epsilon}\boldsymbol{\xi} : \boldsymbol{\xi} \in \{-1, 1\}^d\},\$$

it follows from Lemma 3.25 that f is bounded from above on  $x + \epsilon B$ . In particular, Theorem 3.41 implies that f is Lipschitz continuous on  $x + \delta B$ , provided that  $0 < \delta < \epsilon$ .

Now let K be a compact subset of interior(dom(f)). Then there exists an  $\epsilon > 0$  such that  $K + \epsilon B$  is contained in interior(dom(f)) as well. Since the latter set  $K + \epsilon B$  is compact, too, it follows from continuity of f on interior(dom(f)) that

$$\tilde{L} := \sup_{\boldsymbol{y}, \boldsymbol{z} \in \boldsymbol{K} + \epsilon \boldsymbol{B}} \left| f(\boldsymbol{y}) - f(\boldsymbol{z}) \right| < \infty.$$

Consequently, for arbitrary  $x \in K$  and unit vectors  $v \in B$ ,

$$Df(\boldsymbol{x}, \boldsymbol{v}) \leq \epsilon^{-1} (f(\boldsymbol{x} + \epsilon \boldsymbol{v}) - f(\boldsymbol{x})) \leq \epsilon^{-1} \tilde{L}.$$

As in the proof of Theorem 3.41 one can deduce from this inequality that f is Lipschitz continuous on K with Lipschitz constant  $\epsilon^{-1}\tilde{L}$ .

#### **3.4.2** Subdifferentials and smoothness

A linear functional  $L: \mathbf{V} \to \mathbb{R}$  is called subdifferential of f at the point  $\mathbf{x} \in \text{dom}(f)$  if

(3.11) 
$$f(\boldsymbol{x}) + L(\boldsymbol{v}) \leq f(\boldsymbol{x} + \boldsymbol{v}) \quad \text{for all } \boldsymbol{v} \in \boldsymbol{V}.$$

It follows from our definition of directional derivatives and inequality (3.7) that condition (3.11) is equivalent to

$$(3.12) L(\boldsymbol{v}) \leq Df(\boldsymbol{x}, \boldsymbol{v}) \text{ for all } \boldsymbol{v} \in \boldsymbol{V}.$$

In case of  $Df(x, -v) = -Df(x, v) \in \mathbb{R}$  for arbitrary  $v \in V$ , the function  $Df(x, \cdot)$  itself is the unique subdifferential of f at x.

**Exercise 3.43** (Sublinearity and linearity). Let  $L : V \to \mathbb{R}$  be sublinear, i.e. for arbitrary  $v, w \in V$  and  $\lambda \ge 0$ ,

$$L(\lambda v) = \lambda L(v)$$
 and  $L(v + w) \leq L(v) + L(w)$ .
(a) Show that L is linear if and only if

$$L(-\boldsymbol{v}) = -L(\boldsymbol{v}) \text{ for all } \boldsymbol{v} \in \boldsymbol{V}.$$

(b) Let  $V = \operatorname{span}(b_1, \ldots, b_d)$ . Show that L is linear if and only if

$$L(-\boldsymbol{b}_i) = -L(\boldsymbol{b}_i) \text{ for } 1 \le i \le d$$

**Exercise 3.44** (Sublinearity and Lipschitz-continuity). Let  $L : \mathbb{R}^d \to \mathbb{R}$  be sublinear. Show that for arbitrary  $x, y \in \mathbb{R}^d$ ,

$$|L(\boldsymbol{x}) - L(\boldsymbol{y})| \leq \begin{cases} \max\{L(\boldsymbol{v}) : \boldsymbol{v} \in \{-1, 1\}^d\} \|\boldsymbol{x} - \boldsymbol{y}\|_{\infty}, \\ \max\{L(\boldsymbol{v}) : \boldsymbol{v} \in \{\pm \boldsymbol{e}_1, \dots, \pm \boldsymbol{e}_d\}\} \|\boldsymbol{x} - \boldsymbol{y}\|_1, \end{cases}$$

where  $\|\boldsymbol{x}\|_{\infty} := \max_{i=1,\dots,d} |x_i|, \|\boldsymbol{x}\|_1 := \sum_{i=1}^d |x_i|, \text{ and } \boldsymbol{e}_1, \dots, \boldsymbol{e}_d \text{ is the standard basis of } \mathbb{R}^d.$ 

The next theorem shows that a convex function f on  $\mathbb{R}^d$  is differentiable almost everywhere in the interior of dom(f).

**Theorem 3.45.** Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be convex. For any interior point of dom(f) there exists a subdifferential of f. The set of all interior points of dom(f) at which f is not differentiable is a Borel set with Lebesgue measure 0.

**Proof of Theorem 3.45.** Let  $x \in \Omega$  := interior(dom(f)). We have to show that there exists a linear function L on  $\mathbb{R}^d$  such that  $f(x) + L \leq f(x + \cdot)$ . For this we apply Theorem 2.22 (b) to the disjoint convex subsets  $D_1 := \{x\} \times (-\infty, f(x))$  and  $D_2 := \operatorname{epi}(f)$  of  $\mathbb{R}^d \times \mathbb{R}$ . There exists a nonzero vector  $(w, t) \in \mathbb{R}^d \times \mathbb{R}$  such that

$$\langle (\boldsymbol{x},r), (\boldsymbol{w},t) \rangle \leq \langle (\boldsymbol{x}+\boldsymbol{v},s), (\boldsymbol{w},t) \rangle$$

whenever  $r < f(\boldsymbol{x}), \boldsymbol{v} \in \mathbb{R}^d$  and  $s \geq f(\boldsymbol{x} + \boldsymbol{v})$ . Writing  $\langle (\boldsymbol{x}, r), (\boldsymbol{w}, t) \rangle = \boldsymbol{x}^\top \boldsymbol{w} + rt$  and  $\langle (\boldsymbol{x} + \boldsymbol{v}, s), (\boldsymbol{w}, t) \rangle = (\boldsymbol{x} + \boldsymbol{v})^\top \boldsymbol{w} + st$ , we get the inequalities

(3.13) 
$$rt \leq \mathbf{v}^{\top} \mathbf{w} + st$$
 whenever  $r < f(\mathbf{x}), \mathbf{v} \in \mathbb{R}^d, s \geq f(\mathbf{x} + \mathbf{v}).$ 

For v = 0 we see that  $rt \le st$  for arbitrary  $r < f(x) \le s$ , so  $t \ge 0$ . If t = 0, then  $w \ne 0$ , and (3.13) could be rephrased as  $v^{\top}w \ge 0$  for arbitrary  $v \in \mathbb{R}^d$  such that  $x + v \in \text{dom}(f)$ . But this would contradict our assumption that x is an interior point of dom(f). Thus t > 0, and we may assume without loss of generality that t = 1. Then (3.13) reads

$$f(\boldsymbol{x}) \leq \boldsymbol{v}^{\top} \boldsymbol{w} + f(\boldsymbol{x} + \boldsymbol{v}) \text{ for arbitrary } \boldsymbol{v} \in \mathbb{R}^d,$$

so  $L(v) := -v^{\top}w$  defines a subdifferential of f at x.

Suppose that

$$(3.14) Df(\boldsymbol{x},-\boldsymbol{e}_i) = -Df(\boldsymbol{x},\boldsymbol{e}_i) \text{ for } 1 \le i \le d,$$

where  $e_1, e_2, \ldots, e_d$  is the standard basis of  $\mathbb{R}^d$ . Then  $Df(x, \cdot)$  is linear, see Exercise 3.43. But this implies that f is differentiable at x, see Exercise 3.46. On the other hand, differentiability of f at x implies (3.14).

Consequently, the set of all  $x \in \Omega$  at which f is not differentiable coincides with  $\bigcup_{i=1}^{d} \Omega_i$  with

$$\Omega_i := \{ \boldsymbol{x} \in \Omega : Df(\boldsymbol{x}, -\boldsymbol{e}_i) > -Df(\boldsymbol{x}, \boldsymbol{e}_i) \}.$$

One can easily verify that  $\Omega_i$  is a Borel set, so it remains to be shown that  $\Omega_i$  has Lebesgue measure 0. By Fubini's theorem,

$$\operatorname{Leb}_{d}(\Omega_{i}) = \int_{\boldsymbol{e}_{i}^{\perp}} \operatorname{Leb}_{1}(N_{i}(\boldsymbol{z})) \operatorname{Leb}_{d-1}(d\boldsymbol{z})$$

with

$$N_i(\boldsymbol{z}) := \{ t \in \mathbb{R} : \boldsymbol{z} + t \boldsymbol{e}_i \in \Omega, Df(\boldsymbol{z} + t \boldsymbol{e}_i, -\boldsymbol{e}_i) > -Df(\boldsymbol{z} + t \boldsymbol{e}_i, \boldsymbol{e}_i) \}.$$

But  $g_{z}(t) := f(z + te_i)$  defines a convex function  $g_{z}$  on  $\mathbb{R}$ , and for  $t \in \text{dom}(g_{z})$ ,

$$g'_{\boldsymbol{z}}(t-) = -Df(\boldsymbol{z}+t\boldsymbol{e}_i,-\boldsymbol{e}_i), \quad g'_{\boldsymbol{z}}(t+) = Df(\boldsymbol{z}+t\boldsymbol{e}_i,\boldsymbol{e}_i).$$

As shown in Section 3.1, the set  $N_i(z) = \{t \in \text{dom}(g_z) : g'_z(t-) < gz'(t+)\}$  is at most countable and thus has Lebesgue measure 0.

**Exercise 3.46** ("Subdifferentiability"). Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be convex, and let x be an interior point of dom(f). Show that

$$\sup_{\boldsymbol{v}\in\boldsymbol{B}} \left| \frac{f(\boldsymbol{x}+t\boldsymbol{v}) - f(\boldsymbol{x})}{t} - Df(\boldsymbol{x},\boldsymbol{v}) \right| \to 0 \quad \text{as } t \downarrow 0.$$

## 3.4.3 Jensen's inequality and Bregman divergences

In what follows let  $\mathbf{Y} = (Y_i)_{i=1}^d$  be a random vector such that  $\mathbb{E} \|\mathbf{Y}\| < \infty$ , and let  $\boldsymbol{\mu} := (\mathbb{E} Y_i)_{i=1}^d$  be its mean.

**Theorem 3.47** (Jensen). Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be a convex and measurable function such that  $\mathbb{P}(\mathbf{Y} \in \text{dom}(f)) = 1$ . Then  $\boldsymbol{\mu} \in \text{dom}(f)$  as well, and

$$\mathbb{E} f(\boldsymbol{Y}) \geq f(\boldsymbol{\mu}).$$

If f is strictly convex on dom(f) and  $\mathbb{P}(\mathbf{Y} \neq \boldsymbol{\mu}) > 0$ , then even

$$\mathbb{E} f(\boldsymbol{Y}) > f(\boldsymbol{\mu}).$$

Before proving this theorem, let us verify an auxiliary result.

**Lemma 3.48**. Let  $C \subset \mathbb{R}^d$  be a convex set.

(a) Either C is contained in a hyperplane or it has nonvoid interior. In the latter case,

$$C \subset \overline{\operatorname{interior}(C)} = \overline{C}.$$

(b) Suppose that C is a Borel set with  $\mathbb{P}(Y \in C) = 1$ . Then either

$$\mu \in \operatorname{interior}(C),$$

or there exists a unit vector  $u \in \mathbb{R}^d$  such that  $u^{ op} Y = u^{ op} \mu$  almost surely.

**Proof of Lemma 3.48.** Concerning part (a), let  $x_0$  be an arbitrary point in C, and let

$$oldsymbol{W} \,:=\, \mathrm{span}\{oldsymbol{x} - oldsymbol{x}_0: oldsymbol{x} \in oldsymbol{C}\}$$

If dim(W) < d, there exists a unit vector  $u \in \mathbb{R}^d$  which is perpendicular to W, and  $C \subset x_0 + W$  is contained in the hyperplane  $H := \{x \in \mathbb{R}^d : u^{\top}x = u^{\top}x_0\}$ . If dim(W) = d, the set C has nonvoid interior as shown in Lemma 2.41.

Now we change  $x_0$ , if necessary, such that it is an interior point of C. That means, for some  $\epsilon > 0$ ,

$$x_0 + \epsilon B \subset C.$$

For any point  $x_1 \in C$  let  $x_{\lambda} := (1 - \lambda)x_0 + \lambda x_1$  for  $\lambda \in [0, 1)$ . Then  $x_{\lambda} \to x_1$  as  $\lambda \uparrow 1$ . But  $x_{\lambda}$  is an interior point of C for any  $\lambda \in [0, 1)$ , because

$$oldsymbol{x}_\lambda + (1-\lambda)\epsilonoldsymbol{B} \;=\; (1-\lambda)(oldsymbol{x}_0+\epsilonoldsymbol{B}) + \lambdaoldsymbol{x}_1 \;\subset\; oldsymbol{C}.$$

Thus  $x_1 \in \overline{\operatorname{interior}(C)}$ . This shows that  $C \subset \overline{\operatorname{interior}(C)}$ . The latter inclusion implies that even  $\overline{C} \subset \overline{\operatorname{interior}(C)}$ . Since obviously,  $\overline{\operatorname{interior}(C)} \subset \overline{C}$ , we end up with the equation  $\overline{C} = \overline{\operatorname{interior}(C)}$ .

Now we prove part (b). If C has no interior points, part (a) shows that it is contained in a hyperplane. That means, for some unit vector  $u \in \mathbb{R}^d$  and  $r \in \mathbb{R}$ ,  $C \subset \{x \in \mathbb{R}^d : u^\top x = r\}$ . This implies that  $u^\top Y = r$  almost surely, and taking the expectation of this equation yields  $r = u^\top \mu$ . Now let C have interior points. If  $\mu$  is not an interior point of C, it lies on the boundary or in the complement of  $\overline{C}$ . Since  $\overline{C}$  is a closed and convex set, there exists a unit vector  $u \in \mathbb{R}^d$  such that

$$\boldsymbol{u}^{\top}\boldsymbol{\mu} \geq \boldsymbol{u}^{\top}\boldsymbol{x} \quad ext{for all } \boldsymbol{x} \in \overline{\boldsymbol{C}},$$

see Corollaries 2.19 and 2.21. In particular,  $u^{\top}Y \leq u^{\top}\mu = \mathbb{E}(u^{\top}Y)$  almost surely, and this implies that even  $u^{\top}Y = u^{\top}\mu$  almost surely.

**Proof of Theorem 3.47.** We may assume without loss of generality that  $\mu$  is an interior point of dom(f). Otherwise we could deduce from Lemma 3.48 that  $\mathbb{P}(\mathbf{Y} \in \mathbf{H}) = 1$  for some hyperplane  $\mathbf{H} \subset \mathbb{R}^d$ , so we could replace  $\mathbb{R}^d$  with  $\mathbf{H}$ , and the latter set could be identified with  $\mathbb{R}^{d-1}$ . This reduction could be repeated until  $\mu$  is an interior point of dom(f) and still  $d \ge 1$ , or  $\mathbb{P}(\mathbf{Y} = \mathbf{y}_o) = 1$  for a fixed point  $\mathbf{y}_o \in \mathbb{R}^d$ . The latter case is obviously trivial.

If  $\mu$  is an interior point of dom(f), by Theorem 3.45 there exists a linear function  $L : \mathbb{R}^d \to \mathbb{R}$  such that

$$f(\boldsymbol{y}) \geq f(\boldsymbol{\mu}) + L(\boldsymbol{y} - \boldsymbol{\mu})$$
 for arbitrary  $\boldsymbol{y} \in \mathbb{R}^d$ .

Setting y = Y and taking the expectation of both sides yields that

$$\mathbb{E} f(\mathbf{Y}) \geq f(\boldsymbol{\mu}) + \mathbb{E} L(\mathbf{Y} - \boldsymbol{\mu}) = f(\boldsymbol{\mu}) + L(\underbrace{\mathbb{E}(\mathbf{Y} - \boldsymbol{\mu})}_{=\mathbf{0}}) = f(\boldsymbol{\mu}).$$

If f is strictly convex on dom(f), then even  $f(y) > f(\mu) + L(y - \mu)$  for all  $y \in \mathbb{R}^d \setminus \{\mu\}$ , and we obtain the strict inequality  $\mathbb{E} f(Y) > f(\mu)$ , unless  $\mathbb{P}(Y \neq \mu) = 0$ .  $\Box$ 

The difference  $\mathbb{E} f(\mathbf{Y}) - f(\boldsymbol{\mu})$  can be quantified in the special setting when dom(f) has nonvoid interior and f is continuously differentiable and strictly convex on dom(f). For  $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$  we define

$$D_f(\boldsymbol{y}, \boldsymbol{x}) = f(\boldsymbol{y}) - f(\boldsymbol{x}) - \nabla f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}),$$

the so-called *Bregman divergence* of  $\boldsymbol{y}$  from  $\boldsymbol{x}$  (with respect to f). Note that  $h(t) := f(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x}))$  is strictly convex in  $t \in [0, 1]$  with derivative  $h'(t) = \nabla f(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x}))^{\top} (\boldsymbol{y} - \boldsymbol{x})$ , so

$$D_f(\boldsymbol{y}, \boldsymbol{x}) = h(1) - h(0) - h'(0) = \int_0^1 (h'(t) - h'(0)) dt \ge 0,$$

with equality if and only if y = x.

**Theorem 3.49** (Minimizing expected Bregman divergence). Suppose that  $\mathbb{E} f(\mathbf{Y}) < \infty$ . Then  $\mathbb{E} D_f(\mathbf{Y}, \mathbf{x}) < \infty$  for arbitrary  $\mathbf{x} \in \text{dom}(f)$ , and

$$\mathbb{E} D_f(\boldsymbol{Y}, \boldsymbol{x}) \geq \mathbb{E} D_f(\boldsymbol{Y}, \boldsymbol{\mu})$$

with equality if and only if  $x = \mu$ .

More generally, for arbitrary points  $x, x_o \in \text{dom}(f)$ ,

$$\mathbb{E}(D_f(\boldsymbol{Y}, \boldsymbol{x}) - D_f(\boldsymbol{Y}, \boldsymbol{x}_o))$$

exists in  $\mathbb{R}$  and is minimal with respect to x if and only if  $x = \mu$ .

**Proof of Theorem 3.49.** Recall from Theorem 3.47 that  $\mu \in \text{dom}(f)$ . Note that

$$D_f(\boldsymbol{Y}, \boldsymbol{x}) - D_f(\boldsymbol{Y}, \boldsymbol{x}_o) = f(\boldsymbol{x}_o) + \nabla f(\boldsymbol{x}_o)^\top (\boldsymbol{Y} - \boldsymbol{x}_o) - f(\boldsymbol{x}) - \nabla f(\boldsymbol{x})^\top (\boldsymbol{Y} - \boldsymbol{x})$$

is integrable, because  $\mathbb{E} \| \mathbf{Y} \|$  is finite. The function  $h(\mathbf{x}, \mathbf{x}_o) := \mathbb{E} (D_f(\mathbf{Y}, \mathbf{x}) - D_f(\mathbf{Y}, \mathbf{x}_o))$ satisfies

$$\begin{split} h(\boldsymbol{x}, \boldsymbol{x}_o) - h(\boldsymbol{\mu}, \boldsymbol{x}_o) &= h(\boldsymbol{x}, \boldsymbol{\mu}) \\ &= \mathbb{E} \big[ f(\boldsymbol{\mu}) + \nabla f(\boldsymbol{\mu})^\top (\boldsymbol{Y} - \boldsymbol{\mu}) - f(\boldsymbol{x}) - \nabla f(\boldsymbol{x})^\top (\boldsymbol{Y} - \boldsymbol{x}) \big] \\ &= f(\boldsymbol{\mu}) - f(\boldsymbol{x}) - \nabla f(\boldsymbol{x})^\top (\boldsymbol{\mu} - \boldsymbol{x}) \\ &= D_f(\boldsymbol{\mu}, \boldsymbol{x}). \end{split}$$

This shows that  $\mu$  is the unique minimizer of  $h(\cdot, \boldsymbol{x}_o)$  on dom(f).

Remark 3.50. Note that under the assumptions of the previous theorem,

$$\mathbb{E} f(\boldsymbol{Y}) - f(\boldsymbol{\mu}) = \mathbb{E} \big[ f(\boldsymbol{Y}) - f(\boldsymbol{\mu}) - \nabla f(\boldsymbol{\mu})^{\top} (\boldsymbol{Y} - \boldsymbol{\mu}) \big] = \mathbb{E} D_f(\boldsymbol{Y}, \boldsymbol{\mu}).$$

**Exercise 3.51**. For  $x \in \mathbb{R}^d$  let

$$f(x) := \frac{\|x\|^2}{\sqrt{1+\|x\|^2}}$$

Show that  $f : \mathbb{R}^d \to \mathbb{R}$  is strictly convex and continuously differentiable, and determine the Bregman divergence  $D_f(y, x)$ .

## 3.5 Lower Semicontinuity, Minimizers and Convex Conjugates

So far we know that convex functions on  $\mathbb{R}^d$  behave quite nicely on the interior of their domain. From that perspective it looks a bit strange that we talk about a convex and measurable function in Jensen's inequality. To see that strange things may happen on the boundary, note that *any* function  $f : \mathbb{R}^d \to (-\infty, \infty]$  with

$$f(\boldsymbol{x}) \begin{cases} = 0 & \text{if } \|\boldsymbol{x}\| < 1 \\ \ge 0 & \text{if } \|\boldsymbol{x}\| = 1 \\ = \infty & \text{if } \|\boldsymbol{x}\| > 1 \end{cases}$$

is convex. This follows from the fact that for the standard Euclidean norm  $\|\cdot\|$ , any point on the unit sphere  $\partial B$  is an extremal point of the closed unit ball B. In the next subsection we introduce an additional regularity condition which has some useful consequences if applied to convex functions.

#### 3.5.1 Lower semicontinuity

In this subsection we consider an arbitrary metric space  $(\mathcal{X}, d)$  and study a special property of functions on  $\mathcal{X}$ .

**Definition 3.52** (Lower semicontinuity). A function  $f : \mathcal{X} \to [-\infty, \infty]$  is called *lower semicontinuous* if

$$f(x) \leq \liminf_{y \to x} f(y)$$

for arbitrary  $x \in \mathcal{X}$ .

The previous limes inferior is defined as follows:

$$\liminf_{y \to x} f(y) = \lim_{\delta \to 0+} \inf_{y \in \mathcal{X} : d(x,y) \le \delta} f(y).$$

Note that the limit on the right-hand side exists in  $[-\infty, \infty]$ , because  $\inf_{y: d(x,y) \le \delta} f(y)$  is antitonic in  $\delta > 0$ . Note that for any sequence  $(x_n)_n$  in  $\mathcal{X}$  with limit x,

$$\liminf_{n \to \infty} f(x_n) \ge \liminf_{y \to x} f(y).$$

If we choose  $x_n \in \mathcal{X}$  such that  $d(x, x_n) \leq 1/n$  and  $f(x_n) \leq \inf_{y: d(x,y) \leq 1/n} f(y) + 1/n$ , then

$$\lim_{n \to \infty} f(x_n) = \liminf_{y \to x} f(y).$$

**Lemma 3.53** (Characterizing lower semicontinuity). For a function  $f : \mathcal{X} \to [-\infty, \infty]$  on a metric space  $(\mathcal{X}, d)$  the following three statements are equivalent:

- (i) f is lower semicontinuous.
- (ii) For any  $r \in \mathbb{R}$ , the set  $\{x \in \mathcal{X} : f(x) \leq r\}$  is a closed subset of  $\mathcal{X}$ .
- (iii) The epigraph  $epi(f) = \{(x, r) \in \mathcal{X} \times \mathbb{R} : f(x) \le r\}$  is a closed subset of  $\mathcal{X} \times \mathbb{R}$ .

In part (iii), we equip the cartesian product  $\mathcal{X} \times \mathbb{R}$  with the metric

$$d((x,r),(y,s)) := \sqrt{d(x,y)^2 + |r-s|^2}$$

Alternatively, one could define d((x,r), (y,s)) to be d(x,y) + |r-s| or  $\max(d(x,y), |r-s|)$ .

**Proof of Lemma 3.53.** Suppose that f is lower semicontinuous. If  $(x_n)_n$  is a sequence in  $\{f \le r\}$  with limit  $x \in \mathcal{X}$ , then  $f(x) \le \liminf_{n\to\infty} f(x_n) \le r$ , so  $x \in \{f \le r\}$ , too. This shows that condition (ii) holds true.

Suppose that f satisfies (ii). Let  $((x_n, r_n))_n$  be a sequence in epi(f) with limit  $(x, r) \in \mathcal{X} \times \mathbb{R}$ . For arbitrary fixed r' > r and sufficiently large n,  $f(x_n) \leq r_n \leq r'$ , so x is a limit point of a sequence in  $\{f \leq r'\}$ . Since the latter set is closed by assumption,  $x \in \{f \leq r'\}$ , so  $f(x) \leq r'$ . Since r' can be chosen arbitrarily close to r, we may conclude that  $f(x) \leq r$ , so  $(x, r) \in epi(f)$ . This proves (iii).

Suppose that f satisfies (iii). For any  $x \in \mathcal{X}$  let  $(x_n)_n$  be a sequence in  $\mathcal{X}$  with limit x such that  $r := \lim_{n \to \infty} f(x_n) = \lim_{t \to \infty} \inf_{y \to x} f(y)$ . We have to show that  $r \ge f(x)$ . This is trivial in case of  $r = \infty$ . Otherwise let r' > r. Then for a suitable  $n_o$ ,  $f(x_n) \le r'$  for all  $n \ge n_o$ . Hence  $(x, r') = \lim_{n \to \infty} (x_n, r')$ , and  $(x_n, r') \in \operatorname{epi}(f)$ , for all  $n \ge n_o$ . By assumption (iii), the point (x, r') belongs to  $\operatorname{epi}(f)$ , too, whence  $f(x) \le r'$ . Since r' can be chosen arbitrarily close to r, this shows that  $f(x) \le r$ , i.e. f is lower semicontinuous.

**Exercise 3.54.** Let  $(f_{\lambda})_{\lambda \in \Lambda}$  be a family of functions  $f_{\lambda} : \mathcal{X} \to [-\infty, \infty]$  on some set  $\mathcal{X}$ , and let  $f := \sup_{\lambda \in \Lambda} f_{\lambda}$  (pointwise).

(a) Show that

$$\operatorname{epi}(f) = \bigcap_{\lambda \in \Lambda} \operatorname{epi}(f_{\lambda}).$$

(b) Now suppose that  $(\mathcal{X}, d)$  is a metric space and that all  $f_{\lambda}, \lambda \in \Lambda$ , are lower semicontinuous. Show that f is lower semicontinuous, too.

**Exercise 3.55.** Let  $(\mathcal{X}, d)$  be a metric space, and let  $f : \mathcal{X} \to [-\infty, \infty]$  be lower semicontinuous. Show that for any compact set  $\mathcal{X}_o \subset \mathcal{X}$ ,  $\arg \min_{\mathcal{X}_o} f$  is a compact (and thus nonvoid) set.

## 3.5.2 Convexity, lower semicontinuity and affine functions

Now we return to functions  $f : \mathbb{R}^d \to (-\infty, \infty]$ . As we shall see here and in other parts, combining convexity and lower semicontinuity is particularly fruitful. Recall that f is convex

and lower semicontinuous if and only if its epigraph

$$\operatorname{epi}(f) = \left\{ (\boldsymbol{x}, r) \in \mathbb{R}^d \times \mathbb{R} : f(\boldsymbol{x}) \le r \right\}$$

is a closed and convex subset of  $\mathbb{R}^d \times \mathbb{R}$ . Then epi(f) is an intersection of closed halfspaces in  $\mathbb{R}^d \times \mathbb{R}$ . As shown in Exercise 3.57, the epigraph of an affine function  $A : \mathbb{R}^d \to \mathbb{R}$  is a closed halfspace in  $\mathbb{R}^d \times \mathbb{R}$ . Thus the next result is rather plausible in view of Corollary 2.19:

**Theorem 3.56** (Convex functions as pointwise suprema of affine functions). For a function  $f : \mathbb{R}^d \to (-\infty, \infty]$  the following two statements are equivalent:

(i) f is convex and lower semicontinuous.

(ii) f is the pointwise supremum of all affine functions  $A : \mathbb{R}^d \to \mathbb{R}$  such that  $A \leq f$  on  $\mathbb{R}^d$ .

Note that  $\sup(\emptyset) := -\infty$ , so part (ii) implies that there is some affine function A with  $A \leq f$ .

**Proof of Theorem 3.56.** Let  $\mathcal{A}$  be the set of all affine functions  $A : \mathbb{R}^d \to \mathbb{R}$  such that  $A \leq f$ , and suppose that  $\mathcal{A}$  is nonvoid. Since every function  $A \in \mathcal{A}$  is convex, real-valued and continuous, it follows from Lemma 3.20 (d) and Exercise 3.54 (b) that the pointwise supremum  $\underline{f} := \sup_{A \in \mathcal{A}} A$ is a convex and lower semicontinuous function  $\underline{f} : \mathbb{R}^d \to (-\infty, \infty]$  such that  $\underline{f} \leq f$ . This shows already that condition (ii), which states that  $f \equiv f$ , implies condition (i).

It remains to show that condition (i) implies condition (ii). For this it suffices to show that for any point  $(\boldsymbol{x}_o, r_o) \in \mathbb{R}^d \times \mathbb{R}$  with  $r_o < f(\boldsymbol{x}_o)$  there exists an  $A \in \mathcal{A}$  such that  $A(\boldsymbol{x}_o) \ge r_o$ . In case of  $f \equiv \infty$  this is trivial; just take the constant function  $A \equiv r_o$ . Hence let dom(f) be nonempty. In case of  $\boldsymbol{x}_o \notin \overline{\operatorname{dom}(f)}$  let  $\boldsymbol{x}_*$  be the metric projection of  $\boldsymbol{x}_o$  onto  $\overline{\operatorname{dom}(f)}$ , and pick some number  $r_* < f(\boldsymbol{x}_*)$ . Otherwise let  $(\boldsymbol{x}_*, r_*) := (\boldsymbol{x}_o, r_o)$ . In both cases,

$$\boldsymbol{D} := \operatorname{conv}\{(\boldsymbol{x}_o, r_o), (\boldsymbol{x}_*, r_*)\}$$

is a compact, convex subset of  $\mathbb{R}^d \times \mathbb{R}$  with

$$\boldsymbol{D} \cap \operatorname{epi}(f) = \emptyset$$

Hence by Theorem 2.22 (a) there exists a nonzero vector  $(\boldsymbol{b},c)\in\mathbb{R}^d imes\mathbb{R}$  such that

(3.15) 
$$\max_{(\boldsymbol{x},r)\in\boldsymbol{D}}(\boldsymbol{b}^{\top}\boldsymbol{x}+cr) < \inf_{(\boldsymbol{y},s)\in \operatorname{epi}(f)}(\boldsymbol{b}^{\top}\boldsymbol{y}+cs).$$

Suppose that  $c \leq 0$ . Then consider a sequence  $(\boldsymbol{y}_n)_n$  in dom(f) with limit  $\boldsymbol{x}_*$ . By lower semicontinuity of f,

$$\limsup_{n\to\infty} \left( \boldsymbol{b}^\top \boldsymbol{y}_n + cf(\boldsymbol{y}_n) \right) \leq \boldsymbol{b}^\top \boldsymbol{x}_* + cf(\boldsymbol{x}_*) \leq \boldsymbol{b}^\top \boldsymbol{x}_* + cr_*,$$

a contradiction to (3.15). Hence c > 0 and we may assume without loss of generality that c = 1. Then (3.15) implies that

$$\boldsymbol{b}^{\top} \boldsymbol{x}_o + r_o < \boldsymbol{b}^{\top} \boldsymbol{y} + f(\boldsymbol{y}) \text{ for all } \boldsymbol{y} \in \mathbb{R}^d.$$

In other words,  $A(\boldsymbol{y}) := \boldsymbol{b}^{\top} \boldsymbol{x}_o + r_o - \boldsymbol{b}^{\top} \boldsymbol{y}$  defines an affine function such that A < f and  $A(\boldsymbol{x}_o) = r_o$ .

**Exercise 3.57** (Epigraphs of affine functions). Show that a function  $f : \mathbb{R}^d \to \mathbb{R}$  is affine if and only if its epigraph is a closed halfspace in  $\mathbb{R}^d \times \mathbb{R}$ .

## 3.5.3 Existence of minimizers and coercivity

Let us define an important property in connection with minimization of functions:

**Definition 3.58** (Coercivity). A function  $f : \mathbb{R}^d \to (-\infty, \infty]$  is called *coercive* if

$$f(\boldsymbol{x}) \to \infty$$
 as  $\|\boldsymbol{x}\| \to \infty$ .

The next lemma shows that lower semicontinuity and coercivity of a function imply the existence of a compact set of minimizers:

**Lemma 3.59** (Existence of minimizers). Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be lower semicontinuous and coercive with dom $(f) \neq \emptyset$ . Then

$$rgmin_{oldsymbol{x}\in\mathbb{R}^d} f(oldsymbol{x}) \ = \ ig\{oldsymbol{x}\in\mathbb{R}^d: f(oldsymbol{x}) = \inf_{oldsymbol{y}\in\mathbb{R}^d} f(oldsymbol{y})ig\}$$

is compact (and thus nonvoid).

**Proof of Lemma 3.59.** Let  $(\boldsymbol{x}_n)_n$  be a sequence in  $\mathbb{R}^d$  such that  $f(\boldsymbol{x}_n)$  converges to  $\inf(f) := \inf_{\boldsymbol{y} \in \mathbb{R}^d} f(\boldsymbol{y})$  as  $n \to \infty$ . It follows from f being coercive that  $(\boldsymbol{x}_n)_n$  is bounded. Thus we may replace  $(\boldsymbol{x}_n)_n$  with a subsequence, if necessary, such that  $(\boldsymbol{x}_n)_n$  has a limit  $\boldsymbol{x} \in \mathbb{R}^d$ . By lower semicontinuity of f,

$$f(\boldsymbol{x}) \leq \lim_{n \to \infty} f(\boldsymbol{x}_n) = \inf(f),$$

whence  $f(x) = \inf(f)$ . This shows that  $\arg\min(f) := \{f = \inf(f)\}$  is nonempty. Since  $\arg\min(f)$  equals  $\{f \le \inf(f)\}$ , lower semicontinuity of f implies that it is closed, while coercivity of f implies that it is bounded. Consequently,  $\arg\min(f)$  is compact.

For convex functions one can give a rather complete picture. First of all, for a convex and lower semicontinuous function, coercivity and compactness of its set of minimizers are equivalent.

**Theorem 3.60** (Convexity and coercivity). Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be convex and lower semicontinuous with dom $(f) \neq \emptyset$ . Then f is coercive if and only if the set

$$rgmin_{oldsymbol{x}\in\mathbb{R}^d} f(oldsymbol{x}) \ = \ \Big\{oldsymbol{x}\in\mathbb{R}^d: f(oldsymbol{x}) = \inf_{oldsymbol{y}\in\mathbb{R}^d} f(oldsymbol{y})\Big\}$$

is compact (and thus nonvoid).

For nondegenerate convex functions which are strictly convex on their domain, it is clear that they have at most one minimizer. Hence Theorem 3.60 yields the following corollaries:

**Corollary 3.61.** Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be lower semicontinuous and convex with dom $(f) \neq \emptyset$ . If f is strictly convex on dom(f), then there exists a unique minimizer of f if and only if f is coercive.

Since a real-valued convex function on  $\mathbb{R}^d$  is automatically continuous, there is a simplified version of the previous corollary for this situation:

**Corollary 3.62**. Let  $f : \mathbb{R}^d \to \mathbb{R}$  be strictly convex. There exists a unique minimizer of f if and only if f is coercive.

**Proof of Theorem 3.60.** We know already from Lemma 3.59 that the set  $\arg\min(f) = \{f \le \inf(f)\}$  is compact if f is coercive. It remains to show that compactness of  $\arg\min(f)$  implies coercivity. To this end let  $x_o \in \arg\min(f)$  and let r > 0 such that  $\arg\min(f)$  is contained in  $x_o + rB$ . Since the sphere with center  $x_o$  and radius r + 1 is compact, lower semicontinuity of f implies the existence of

$$c := \min_{\boldsymbol{x} \in \mathbb{R}^d : \|\boldsymbol{x} - \boldsymbol{x}_o\| = r+1} f(\boldsymbol{x}) > f(\boldsymbol{x}_o);$$

see also Exercise 3.55 or the proof of Lemma 3.59. For any  $x \in \mathbb{R}^d$  with  $||x - x_o|| \ge r + 1$  we may write  $x = x_o + ||x - x_o||u$  with a unit vector u, and by convexity of  $f(x_o + tu)$  in  $t \in \mathbb{R}$ ,

$$egin{aligned} f(m{x}) &= f(m{x}_o) + ig(f(m{x}_o + \|m{x} - m{x}_o\|m{u}) - f(m{x}_o)ig) \ &\geq f(m{x}_o) + \|m{x} - m{x}_o\|rac{f(m{x}_o + (r+1)m{u}) - f(m{x}_o)}{r+1} \ &\geq f(m{x}_o) + \|m{x} - m{x}_o\|rac{c-f(m{x}_o)}{r+1} \ &\geq f(m{x}_o) + ig(\|m{x}\| - \|m{x}_o\|ig)rac{c-f(m{x}_o)}{r+1}. \end{aligned}$$

The right-hand side converges to  $\infty$  as  $||x|| \to \infty$ , whence f is coercive.

We end this subsection with a specific characterization of coercivity in terms of directional derivatives which is often rather useful:

**Lemma 3.63** (Criterion for coercivity). Let  $f : \mathbb{R}^d \to \mathbb{R}$  be convex. Then f is coercive if and only if for any fixed unit vector  $u \in \mathbb{R}^d$ ,

$$(3.16) \qquad \qquad \lim_{t\to\infty} Df(t\boldsymbol{u},\boldsymbol{u}) > 0.$$

Note that convexity of f implies that for a fixed unit vector  $\boldsymbol{u} \in \mathbb{R}^d$ , the function  $h_{\boldsymbol{u}} : \mathbb{R} \to \mathbb{R}$ ,  $h_{\boldsymbol{u}}(t) := f(t\boldsymbol{u})$ , is convex with  $h'_{\boldsymbol{u}}(t+) = Df(t\boldsymbol{u}, \boldsymbol{u})$ . Thus  $Df(t\boldsymbol{u}, \boldsymbol{u})$  is isotonic in t, and the limit  $\lim_{t\to\infty} Df(t\boldsymbol{u}, \boldsymbol{u})$  exists in  $(-\infty, \infty]$ .

**Proof of Lemma 3.63.** Suppose first that f is coercive. Then for sufficiently large r > 0,

$$c := \min_{\boldsymbol{x} \in r \partial \boldsymbol{B}} f(\boldsymbol{x}) > f(\boldsymbol{0}).$$

But by convexity of f, for any fixed unit vector  $u \in \partial B$ , the function  $h_u(t) := f(tu)$  is convex in  $t \in \mathbb{R}$  with  $h'_u(t+) = Df(tu, u)$  being isotonic in t. Thus

$$\lim_{t \to \infty} Df(t\boldsymbol{u}, \boldsymbol{u}) \geq h'_{\boldsymbol{u}}(r+) \geq r^{-1}(h_{\boldsymbol{u}}(r) - h_{\boldsymbol{u}}(0)) \geq r^{-1}(c - f(\boldsymbol{0})) > 0.$$

Hence (3.16) is satisfied.

Now suppose that f fails to be coercive. That means, for some sequence  $(x_n)_n$  in  $\mathbb{R}^d \setminus \{0\}$  and some real constant C,

$$\lim_{n \to \infty} \|\boldsymbol{x}_n\| = \infty \quad \text{but} \quad f(\boldsymbol{x}_n) \leq C \text{ for all } n.$$

We may replace  $(\boldsymbol{x}_n)_n$  with a subsequence, if necessary, such that in addition

$$oldsymbol{u}_n := \|oldsymbol{x}_n\|^{-1}oldsymbol{x}_n \ o \ oldsymbol{u} \in \partial oldsymbol{B} \ ext{ as } n o \infty.$$

But then for any fixed t > 0, convexity and continuity of f imply that

$$Df(t\boldsymbol{u},\boldsymbol{u}) \leq f((t+1)\boldsymbol{u}) - f(t\boldsymbol{u})$$

$$= \lim_{n \to \infty} \left( f((t+1)\boldsymbol{u}_n) - f(t\boldsymbol{u}_n) \right)$$

$$\leq \lim_{n \to \infty} \frac{f(\|\boldsymbol{x}_n\| \boldsymbol{u}_n) - f(t\boldsymbol{u}_n)}{\|\boldsymbol{x}_n\| - t}$$

$$= \lim_{n \to \infty} \frac{f(\boldsymbol{x}_n) - f(t\boldsymbol{u}_n)}{\|\boldsymbol{x}_n\| - t}$$

$$\leq \lim_{n \to \infty} \frac{C - \min_{\boldsymbol{y} \in t\boldsymbol{B}} f(\boldsymbol{y})}{\|\boldsymbol{x}_n\| - t}$$

$$= 0.$$

Thus (3.16) is violated as well.

**Example 3.64** (Logistic regression). In logistic regression one considers observations  $(z_1, y_1)$ , ...,  $(z_n, y_n)$  in  $\mathbb{R}^d \times \{0, 1\}$  and wants to minimize the negative log-likelihood function  $f : \mathbb{R}^p \to \mathbb{R}$  given by

$$f(\boldsymbol{x}) := \sum_{i=1}^n \left( -y_i \boldsymbol{z}_i^\top \boldsymbol{x} + \log(1 + e^{\boldsymbol{z}_i^\top \boldsymbol{x}}) \right).$$

An obvious question is whether there exists a unique minimizer  $x_*$  of f.

Fist of all, elementary calculations reveal that the gradient and Hessian matrix of f are given by

(3.17) 
$$\nabla f(\boldsymbol{x}) = \sum_{i=1}^{n} (\ell(\boldsymbol{z}_{i}^{\top}\boldsymbol{x}) - y_{i}) \boldsymbol{z}_{i},$$

(3.18) 
$$D^2 f(\boldsymbol{x}) = \sum_{i=1}^n \ell'(\boldsymbol{z}_i^\top \boldsymbol{x}) \, \boldsymbol{z}_i \boldsymbol{z}_i^\top,$$

where  $\ell : \mathbb{R} \to (0, 1)$  is the logistic distribution function,

$$\ell(t) := \frac{e^t}{1+e^t} = (1+e^{-t})^{-1},$$

with derivative  $\ell' = \ell(1 - \ell) > 0$ . Hence for any  $\boldsymbol{v} \in \mathbb{R}^p$ ,

This shows that f is convex, and it is even strictly convex if

$$\operatorname{span}(\boldsymbol{z}_1,\ldots,\boldsymbol{z}_n) = \mathbb{R}^d.$$

As to coercivity, for any fixed unit vector  $\boldsymbol{u} \in \mathbb{R}^d$  and  $t \in \mathbb{R}$ ,

$$Df(t\boldsymbol{u},\boldsymbol{u}) = \nabla f(t\boldsymbol{u})^{\top}\boldsymbol{u}$$
  
=  $\sum_{i=1}^{n} (\ell(t\boldsymbol{z}_{i}^{\top}\boldsymbol{u}) - y_{i})\boldsymbol{z}_{i}^{\top}\boldsymbol{u}$   
 $\rightarrow \sum_{i=1}^{n} ((1 - y_{i})(\boldsymbol{z}_{i}^{\top}\boldsymbol{u})^{+} + y_{i}(\boldsymbol{z}_{i}^{\top}\boldsymbol{u})^{-})$  as  $t \rightarrow \infty$ .

The latter limit is nonnegative, and it is strictly positive unless

$$y_i = \begin{cases} 1 & \text{if } \boldsymbol{z}_i^\top \boldsymbol{u} > 0, \\ 0 & \text{if } \boldsymbol{z}_i^\top \boldsymbol{u} < 0. \end{cases}$$

Hence the function f is coercive if and only if for each unit vector  $\boldsymbol{u} \in \mathbb{R}^d$  there exists at least one index i such that  $y_i = 0 < \boldsymbol{z}_i^\top \boldsymbol{u}$  or  $1 - y_i = 0 > \boldsymbol{z}_i^\top \boldsymbol{u}$ .

In other words, the function f fails to be coercive if and only if there exists a unit vector  $u \in \mathbb{R}^d$  such that

$$\{\boldsymbol{z}_i: y_i = 0\} \subset \{\boldsymbol{x}: \boldsymbol{u}^\top \boldsymbol{x} \le 0\} \text{ and } \{\boldsymbol{z}_i: y_i = 1\} \subset \{\boldsymbol{x}: \boldsymbol{u}^\top \boldsymbol{x} \ge 0\}.$$

Exercise 3.65 (Logistic regression, I). Verify formulae (3.17) and (3.18).

**Exercise 3.66** (Least absolute deviations regression). For given observations  $(z_1, y_1)$ ,  $(z_2, y_2)$ , ...,  $(z_n, y_n)$  in  $\mathbb{R}^d \times \mathbb{R}$  let  $f : \mathbb{R}^d \to \mathbb{R}$  be given by

$$f(oldsymbol{x}) \ := \ \sum_{i=1}^n |oldsymbol{z}_i^ op oldsymbol{x} - y_i|$$

- (a) Show that f is convex.
- (b) Show that f is coercive if and only if

$$\operatorname{span}(\boldsymbol{z}_1,\ldots,\boldsymbol{z}_n) = \mathbb{R}^d.$$

## 3.5.4 Convex conjugates

**Definition 3.67** (Fenchel–Lagrange transform). Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  with dom $(f) \neq \emptyset$ . The Fenchel–Legendre transform of f (also called *convex conjugate of* f) is defined as the function  $f^* : \mathbb{R}^d \to (-\infty, \infty]$ ,

$$f^*(oldsymbol{y}) \ \coloneqq \ \sup_{oldsymbol{x} \in \mathbb{R}^d} (oldsymbol{x}^ op oldsymbol{y} - f(oldsymbol{x})) \ = \ \sup_{oldsymbol{x} \in ext{dom}(f)} (oldsymbol{x}^ op oldsymbol{y} - f(oldsymbol{x})).$$

It follows from Lemma 3.20 (c) and Exercise 3.54 that the convex conjugate  $f^*$  of f is convex and lower semicontinuous. Moreover, the definition of  $f^*$  implies that

Now consider an affine function  $A : \mathbb{R}^d \to \mathbb{R}$ ,  $A(\boldsymbol{x}) := \boldsymbol{b}^\top \boldsymbol{x} - c$ , with a certain vector  $\boldsymbol{b} \in \mathbb{R}^d$ and a real number c. Then  $A \leq f$  if and only if

$$\boldsymbol{x}^{\top}\boldsymbol{b} - f(\boldsymbol{x}) \leq c \quad \text{for all } \boldsymbol{x} \in \mathbb{R}^d,$$

and this is equivalent to

$$f^*(\boldsymbol{b}) \leq c.$$

Consequently, if  $\mathcal{A}$  denotes the set of all affine functions A on  $\mathbb{R}^d$  with  $A \leq f$ , then

$$\underline{f}(\boldsymbol{x}) := \sup_{A \in \mathcal{A}} A(\boldsymbol{x}) = \sup_{(\boldsymbol{b}, c) \in \mathbb{R}^d \times \mathbb{R} : c \ge f^*(\boldsymbol{b})} (\boldsymbol{b}^\top \boldsymbol{x} - c)$$
$$= \sup_{\boldsymbol{b} \in \mathbb{R}^d : f^*(\boldsymbol{b}) < \infty} (\boldsymbol{b}^\top \boldsymbol{x} - f^*(\boldsymbol{b}))$$
$$= (f^*)^*(\boldsymbol{x}).$$

These considerations and Theorem 3.56 yield the following result:

**Theorem 3.68**. Let  $f : \mathbb{R}^d \to (-\infty, \infty]$  be convex and lower semicontinuous with dom $(f) \neq \emptyset$ . Then  $f^*$  has these properties, too, and

$$(f^*)^* \equiv f$$

**Example 3.69**. Let  $f(x) := ||x||^2/2$ . Then

$$x^{\top}y - f(x) = ||y||^2/2 - ||x - y||^2/2 \le f(y)$$

with equality if and only if x = y. Thus

$$f^* \equiv f.$$

Exercise 3.70 (Conjugates of quadratic functions). Let

$$f(\boldsymbol{x}) := c + \boldsymbol{b}^{\top} \boldsymbol{x} + 2^{-1} \boldsymbol{x}^{\top} \boldsymbol{A} \boldsymbol{x}$$

with  $c \in \mathbb{R}$ ,  $b \in \mathbb{R}^d$  and  $A \in \mathbb{R}^{d \times d}$  symmetric and positive definite. Show that

$$f^*(\boldsymbol{y}) = c_* + \boldsymbol{b}_*^\top \boldsymbol{y} + 2^{-1} \boldsymbol{y}^\top \boldsymbol{A}_* \boldsymbol{y}$$

with  $c_* := 2^{-1} \boldsymbol{b}^\top \boldsymbol{A}^{-1} \boldsymbol{b} - c$ ,  $\boldsymbol{b}_* := -\boldsymbol{A}^{-1} \boldsymbol{b}$  and  $\boldsymbol{A}_* := \boldsymbol{A}^{-1}$ .

The following two exercises are about the standard *p*-norms on  $\mathbb{R}^d$ : For a vector  $x \in \mathbb{R}^d$  and  $p \in [1, \infty)$  we write

$$\|\boldsymbol{x}\|_p := \Big(\sum_{i=1}^d |x_i|^p\Big)^{1/p},$$

and

$$\|\boldsymbol{x}\|_{\infty} := \max_{i=1,\dots,d} |x_i|.$$

Exercise 3.71 (Dual norms, I, and Young's inequality). Let

$$f(x) := ||x||_p^p/p$$

with 1 . Show that

$$f^*(\boldsymbol{y}) = \|\boldsymbol{y}\|_q^q/q$$

with q := p/(p-1), so  $p^{-1} + q^{-1} = 1$ . Then show that

$$|oldsymbol{x}^{ op}oldsymbol{y}| \leq f(oldsymbol{x}) + f^*(oldsymbol{y})$$
 and  $|oldsymbol{x}^{ op}oldsymbol{y}| \leq \|oldsymbol{x}\|_p \|oldsymbol{y}\|_q$ 

for arbitrary  $x, y \in \mathbb{R}^d$ .

**Exercise 3.72** (Dual norms, II). Show that for  $p \in [1, \infty]$ , the convex conjugate of the function  $f(x) := \|x\|_p$  is given by

$$f^*(\boldsymbol{y}) = \begin{cases} 0 & \text{if } \|\boldsymbol{y}\|_q \leq 1, \\ \infty & \text{else,} \end{cases}$$

where

$$q := \begin{cases} \infty & \text{if } p = 1, \\ p/(p-1) & \text{if } 1$$

**Remark 3.73** (Support functions). Let C be a nonvoid convex and closed subset of  $\mathbb{R}^d$ , and define

$$f_{\boldsymbol{C}}(\boldsymbol{x}) := \begin{cases} 0 & \text{if } \boldsymbol{x} \in \boldsymbol{C}, \\ \infty & \text{else.} \end{cases}$$

Then

$$f^*_{\boldsymbol{C}}(\boldsymbol{y}) = \sup_{\boldsymbol{x} \in \boldsymbol{C}} \boldsymbol{x}^\top \boldsymbol{y} = h_{\boldsymbol{C}}(\boldsymbol{y}).$$

In case of C being a closed and convex cone,

$$f^*_{\boldsymbol{C}} = f_{\boldsymbol{C}^*}$$

with the polar cone

$$oldsymbol{C}^* \;=\; ig\{oldsymbol{y} \in \mathbb{R}^d: oldsymbol{x}^ op oldsymbol{y} \leq 0 ext{ for all } oldsymbol{x} \in oldsymbol{C}ig\}$$

introduced in Remark 2.26.

Minimizing sums of convex functions. Sometimes one wants to minimize the sum of two convex functions  $f, g : \mathbb{R}^d \to (-\infty, \infty]$ , where we assume that  $\operatorname{dom}(f) \cap \operatorname{dom}(g) \neq \emptyset$ . It may happen that the minimization of  $\boldsymbol{y} \mapsto f^*(\boldsymbol{y}) + g^*(-\boldsymbol{y})$  is easier. Under the slightly stronger constraint that interior $(\operatorname{dom}(f)) \cap \operatorname{dom}(g) \neq \emptyset$ , both minimization problems are strongly related:

**Theorem 3.74** (Fenchel). Let  $f, g : \mathbb{R}^d \to (-\infty, \infty]$  be convex functions such that the intersection of interior  $(\operatorname{dom}(f))$  and  $\operatorname{dom}(g)$  is nonvoid. Then

$$\inf_{oldsymbol{x}\in\mathbb{R}^d}ig(f(oldsymbol{x})+g(oldsymbol{x})ig) \ = \ -\inf_{oldsymbol{y}\in\mathbb{R}^d}ig(f^*(oldsymbol{y})+g^*(-oldsymbol{y})ig).$$

The latter infimum is a minimum if it is real.

A special case of this theorem is the minimization of a convex function f over a closed convex cone C. With its polar cone  $C^*$ ,

$$\inf_{\boldsymbol{x}\in\boldsymbol{C}} f(\boldsymbol{x}) \ = \ -\inf_{\boldsymbol{y}\in-\boldsymbol{C}^*} \ f^*(\boldsymbol{y}).$$

**Proof of Theorem 3.74.** Let  $\alpha := \inf(f + g)$ . Then

$$\begin{aligned} \alpha &= -\sup_{\boldsymbol{x} \in \mathbb{R}^d} \left( \boldsymbol{0}^\top \boldsymbol{x} - f(\boldsymbol{x}) - g(\boldsymbol{x}) \right) \\ &= -\inf_{\boldsymbol{y} \in \mathbb{R}^d} \sup_{\boldsymbol{x} \in \mathbb{R}^d} \left( \boldsymbol{y}^\top \boldsymbol{x} - f(\boldsymbol{x}) + (-\boldsymbol{y})^\top \boldsymbol{x} - g(\boldsymbol{x}) \right) \\ &\geq -\inf_{\boldsymbol{y} \in \mathbb{R}^d} (f^*(\boldsymbol{y}) + g^*(-\boldsymbol{y})). \end{aligned}$$

Hence it suffices to show that

(3.19) 
$$\inf_{\boldsymbol{y}\in\mathbb{R}^d} (f^*(\boldsymbol{y}) + g^*(-\boldsymbol{y})) \leq -\alpha.$$

Since  $\operatorname{interior}(\operatorname{dom}(f)) \cap \operatorname{dom}(g) \neq \emptyset$ , we know that  $\alpha < \infty$ . In case of  $\alpha = -\infty$ , inequality (3.19) is trivial, so it suffices to consider the case  $\alpha \in \mathbb{R}$ . Since  $f + g \ge \alpha$ , the two sets

$$oldsymbol{D}_1 \ := \ ig\{(oldsymbol{x},s):oldsymbol{x}\in\mathbb{R}^d,s>f(oldsymbol{x})ig\}$$

and

$$\boldsymbol{D}_2 := \left\{ (\boldsymbol{y},t) : \boldsymbol{y} \in \mathbb{R}^d, t \leq \alpha - g(\boldsymbol{y}) \right\}$$

are convex and disjoint. Thus there exists a nonzero vector  $(\boldsymbol{b},c)\in\mathbb{R}^d imes\mathbb{R}$  such that

$$\boldsymbol{b}^{ op} \boldsymbol{x} + cs \ \leq \ \boldsymbol{b}^{ op} \boldsymbol{y} + ct \quad ext{for arbitrary} \ (\boldsymbol{x},s) \in \boldsymbol{D}_1, (\boldsymbol{y},t) \in \boldsymbol{D}_2.$$

Taking  $\boldsymbol{x} = \boldsymbol{y} = \boldsymbol{x}_o \in \operatorname{interior}(\operatorname{dom}(f)) \cap \operatorname{dom}(g)$  and  $f(\boldsymbol{x}_o) < s \to \infty$ ,  $\alpha - g(\boldsymbol{x}_o) \ge t \to -\infty$ shows that  $c \le 0$ . If c = 0, we could take  $\boldsymbol{x} = \boldsymbol{x}_o + \delta \boldsymbol{b}$  and  $\boldsymbol{y} = \boldsymbol{x}_o$  with  $\delta > 0$  sufficiently small such that  $\boldsymbol{x} \in \operatorname{dom}(f)$ . This would lead to the contradiction  $\boldsymbol{b} = \boldsymbol{0}$ . Hence c < 0, and we may assume without loss of generality that c = -1. Then we may conclude that

$$\boldsymbol{b}^{\top}\boldsymbol{x} - f(\boldsymbol{x}) \leq \boldsymbol{b}^{\top}\boldsymbol{y} - \alpha + g(\boldsymbol{y}) \text{ for all } \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^{d}$$

Rewriting the latter inequality as  $\boldsymbol{b}^{\top}\boldsymbol{x} - f(\boldsymbol{x}) + (-\boldsymbol{b})^{\top}\boldsymbol{y} - g(\boldsymbol{y}) \leq -\alpha$ , we realize that

$$f^*(oldsymbol{b}) + g^*(-oldsymbol{b}) \ \le \ -lpha \ \le \ \inf_{oldsymbol{y} \in \mathbb{R}^d} ig(f^*(oldsymbol{y}) + g^*(-oldsymbol{y})ig).$$

Consequently, inequality (3.19) is an equality, and  $\boldsymbol{b}$  is a minimizer of the function  $\boldsymbol{y} \mapsto f^*(\boldsymbol{y}) + g^*(-\boldsymbol{y})$ .

## **Chapter 4**

# **Multivariate Optimization**

In this chapter we modify methods from Chapter 1 for multivariate functions.

## 4.1 Newton's Method

Let  $\mathbf{f} = (f_i)_{i=1}^d$  be a continuously differentiable mapping from an open set  $\mathcal{X} \subset \mathbb{R}^d$  to  $\mathbb{R}^d$ . We assume that the Jacobian matrix

$$Doldsymbol{f}(oldsymbol{x}) \, := \, \left(rac{\partial f_i(oldsymbol{x})}{\partial x_j}
ight)_{i,j=1}^d \, \in \, \mathbb{R}^{d imes d}$$

is nonsingular for all  $x \in \mathcal{X}$ . Now we are looking for a zero  $x_*$  of f, that means, a point  $x_* \in f^{-1}(\mathbf{0})$ .

Let  $x_o$  be a first candidate for a zero. We assume that we are already rather close to a zero  $x_*$  and approximate f by an affine function

$$\boldsymbol{x} \mapsto \boldsymbol{f}(\boldsymbol{x}_o) + D\boldsymbol{f}(\boldsymbol{x}_o)(\boldsymbol{x} - \boldsymbol{x}_o).$$

The latter function has a unique zero at

$$\boldsymbol{\psi}(\boldsymbol{x}_o) := \boldsymbol{x}_o - D\boldsymbol{f}(\boldsymbol{x}_o)^{-1}\boldsymbol{f}(\boldsymbol{x}_o),$$

and we hope that  $\psi(x_o)$  is even closer to  $x_*$  than  $x_o$ .

Indeed, iterating this mapping leads to a sequence converging rapidly to a zero of f, provided that we start sufficiently close to the latter. In what follows we work with the Euclidean norm

$$\|\boldsymbol{w}\| := \sqrt{\boldsymbol{w}^\top \boldsymbol{w}},$$

of vectors  $oldsymbol{w} \in \mathbb{R}^{\ell}$ , and for a matrix  $oldsymbol{A} \in \mathbb{R}^{k imes d}$  let

$$\|oldsymbol{A}\| \, := \, \sup_{oldsymbol{v}\in\mathbb{R}^d\,:\,\|oldsymbol{v}\|\leq 1}\,\|oldsymbol{A}oldsymbol{v}\| = \, \sup_{oldsymbol{v}\in\mathbb{R}^d\,:\,\|oldsymbol{v}\|\leq 1,\,oldsymbol{w}\in\mathbb{R}^k\,:\,\|oldsymbol{w}\|\leq 1}\,oldsymbol{w}^ opoldsymbol{A}oldsymbol{v},$$

the so-called operator norm of A. The definition implies that

$$\|Av\| \leq \|A\| \|v\|$$
 for all  $v \in \mathbb{R}^d$ .

An important fact is that any identity matrix has operator norm one, and for arbitrary matrices A, B whose product AB is well-defined,

$$||AB|| \leq ||A|| ||B||.$$

These are the main advantages of the operator norm over the Frobenius norm. But these two matrix norms are equivalent:

**Exercise 4.1.** Show that the operator norm  $||\mathbf{A}||$  and Frobenius norm  $||\mathbf{A}||_F$  of a matrix  $\mathbf{A} \in \mathbb{R}^{k \times d}$  satisfy the inequalities

$$\|\boldsymbol{A}\| \leq \|\boldsymbol{A}\|_F \leq \sqrt{\min(k,d)} \|\boldsymbol{A}\|.$$

**Theorem 4.2** (Local convergence of Newton's method). Let  $x_*$  be a zero of f.

(a) Then

$$\lim_{\boldsymbol{x} \to \boldsymbol{x}_*, \boldsymbol{x} \neq \boldsymbol{x}_*} \frac{\|\boldsymbol{\psi}(\boldsymbol{x}) - \boldsymbol{x}_*\|}{\|\boldsymbol{x} - \boldsymbol{x}_*\|} = 0.$$

(b) Let Df be Lipschitz continuous in a neighborhood U of  $x_*$  with constant L, that means,

$$\|D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{y})\| \leq L\|\boldsymbol{x} - \boldsymbol{y}\|$$
 for all  $\boldsymbol{x}, \boldsymbol{y} \in U_{\boldsymbol{x}}$ 

Then

$$\limsup_{\boldsymbol{x} \to \boldsymbol{x}_*, \boldsymbol{x} \neq \boldsymbol{x}_*} \frac{\|\boldsymbol{\psi}(\boldsymbol{x}) - \boldsymbol{x}_*\|}{\|\boldsymbol{x} - \boldsymbol{x}_*\|^2} \leq \frac{L\|D\boldsymbol{f}(\boldsymbol{x}_*)^{-1}\|}{2}.$$

Consequently, with the multivariate version of Newton's method we achieve superlinear or even quadratic convergence. By means of Exercise 4.1 one can verify that a sufficient condition for Lipschitz-continuity of Df on a convex set  $U \subset \mathcal{X}$  is that all second derivatives  $\partial^2 f_i / (\partial x_k \partial x_\ell)$  are bounded on U. In the proof of Theorem 4.2 we use the following inequalities for matrices:

**Lemma 4.3**. Let A be a nonsingular matrix in  $\mathbb{R}^{d \times d}$ , and let  $B = A + \Delta$  for some  $\Delta \in \mathbb{R}^{d \times d}$ . In case of  $\|\Delta\| < \|A^{-1}\|^{-1}$ , the matrix B is nonsingular, too, and

$$egin{array}{ll} \|m{B}^{-1}\| &\leq & rac{\|m{A}^{-1}\|}{1-\|m{A}^{-1}\|\|\Delta\|}, \ \|m{B}^{-1}-m{A}^{-1}\| &\leq & rac{\|m{A}^{-1}\|^2\|\Delta\|}{1-\|m{A}^{-1}\|\|\Delta\|}. \end{array}$$

We also need an inequality for vector-valued integrals:

**Exercise 4.4.** Let  $g: [a, b] \to \mathbb{R}^d$  have integrable components  $g_i: [a, b] \to \mathbb{R}$ . Show that

$$\left\|\int_{a}^{b} \boldsymbol{g}(t) dt\right\| \leq \int_{a}^{b} \left\|\boldsymbol{g}(t)\right\| dt < \infty,$$

where the integral on the left hand side is defined component-wise.

**Proof of Lemma 4.3.** We rewrite *B* as

$$\boldsymbol{B} = \boldsymbol{A} + \boldsymbol{\Delta} = \boldsymbol{A}(\boldsymbol{I} + \boldsymbol{A}^{-1}\boldsymbol{\Delta}).$$

A matrix of type I + M is nonsingular with inverse

$$(I + M)^{-1} = \sum_{i=0}^{\infty} (-1)^{i} M^{i},$$

provided that ||M|| < 1. In our context, we may apply this to  $M = A^{-1}\Delta$  with norm  $||M|| \le ||A^{-1}|| ||\Delta||$ . Hence B is nonsingular with inverse

$$\boldsymbol{B}^{-1} \;=\; (\boldsymbol{I} + \boldsymbol{A}^{-1} \Delta)^{-1} \boldsymbol{A}^{-1} \;=\; \sum_{i=0}^{\infty} (-1)^{i} (\boldsymbol{A}^{-1} \Delta)^{i} \boldsymbol{A}^{-1} \;=\; \boldsymbol{A}^{-1} + \sum_{i=1}^{\infty} (-1)^{i} (\boldsymbol{A}^{-1} \Delta)^{i} \boldsymbol{A}^{-1},$$

provided that  $\|A^{-1}\|\|\Delta\| < 1$ . The asserted inequalities follow from the fact that

$$\left\|\sum_{i=k}^{\infty} (-1)^{i} (\boldsymbol{A}^{-1} \Delta)^{i} \boldsymbol{A}^{-1}\right\| \leq \sum_{i=k}^{\infty} \|\boldsymbol{A}^{-1}\|^{i+1} \|\Delta\|^{i} = \frac{\|\boldsymbol{A}^{-1}\|^{k+1} \|\Delta\|^{k}}{1 - \|\boldsymbol{A}^{-1}\|\|\Delta\|}$$
  
0, 1.

**Proof of Theorem 4.2.** It follows from continuity of  $x \mapsto Df(x)$  that

$$ho(\delta) \ := \ \sup_{oldsymbol{x}\in\mathcal{X}\,:\,\|oldsymbol{x}-oldsymbol{x}_*\|\leq\delta} \left\|Doldsymbol{f}(oldsymbol{x})-Doldsymbol{f}(oldsymbol{x}_*)
ight\|\ o\ 0 \ \ ext{as}\ \delta\downarrow 0.$$

For sufficiently small  $\delta > 0$ , the closed ball  $B_{\delta}(\boldsymbol{x}_*)$  with center  $\boldsymbol{x}_*$  and radius  $\delta$  is contained in  $\mathcal{X}$ , and for any vector  $\boldsymbol{x} \in B_{\delta}(\boldsymbol{x}_*)$ ,

$$\begin{split} \psi(\boldsymbol{x}) - \boldsymbol{x}_* &= \boldsymbol{x} - \boldsymbol{x}_* - D\boldsymbol{f}(\boldsymbol{x})^{-1}\boldsymbol{f}(\boldsymbol{x}) \\ &= \boldsymbol{x} - \boldsymbol{x}_* - D\boldsymbol{f}(\boldsymbol{x})^{-1} \int_0^1 D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*))(\boldsymbol{x} - \boldsymbol{x}_*) \, dt \\ &= D\boldsymbol{f}(\boldsymbol{x})^{-1} \int_0^1 \left( D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*)) \right) (\boldsymbol{x} - \boldsymbol{x}_*) \, dt. \end{split}$$

This implies that

for k =

$$\begin{aligned} \|\psi(\boldsymbol{x}) - \boldsymbol{x}_*\| &\leq \|D\boldsymbol{f}(\boldsymbol{x})^{-1}\| \left\| \int_0^1 (D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*)))(\boldsymbol{x} - \boldsymbol{x}_*) dt \right\| \\ &\leq \|D\boldsymbol{f}(\boldsymbol{x})^{-1}\| \int_0^1 \| (D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*)))(\boldsymbol{x} - \boldsymbol{x}_*) \| dt \\ &\leq \|D\boldsymbol{f}(\boldsymbol{x})^{-1}\| \int_0^1 \|D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*))\| dt \|\boldsymbol{x} - \boldsymbol{x}_*\|, \end{aligned}$$
(4.1)

where we used the inequality from Exercise 4.4 in the second step. But it follows from Lemma 4.3 that || = c(-) = 1

$$\|D\boldsymbol{f}(\boldsymbol{x})^{-1}\| \leq \frac{\|D\boldsymbol{f}(\boldsymbol{x}_*)^{-1}\|}{1 - \|D\boldsymbol{f}(\boldsymbol{x}_*)^{-1}\|\rho(\delta)},$$

provided that  $\rho(\delta) < \| D \boldsymbol{f}(\boldsymbol{x}_*)^{-1} \|^{-1},$  and

$$\begin{split} \left\| D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*)) \right\| \\ &\leq \left\| D\boldsymbol{f}(\boldsymbol{x}) - D\boldsymbol{f}(\boldsymbol{x}_*) \right\| + \left\| D\boldsymbol{f}(\boldsymbol{x}_* + t(\boldsymbol{x} - \boldsymbol{x}_*)) - D\boldsymbol{f}(\boldsymbol{x}_*) \right\| \\ &\leq 2\rho(\delta). \end{split}$$

Hence for sufficiently small  $\delta > 0$  and  $x \neq x_*$ ,

$$\frac{\|\psi(\bm{x}) - \bm{x}_*\|}{\|\bm{x} - \bm{x}_*\|} \ \le \ \frac{2\|D\bm{f}(\bm{x}_*)^{-1}\|\rho(\delta)}{1 - \|D\bm{f}(\bm{x}_*)^{-1}\|\rho(\delta)}.$$

Since the latter bound tends to 0 as  $\delta \downarrow 0$ , this proves part (a).

For part (b) we choose  $\delta > 0$  small enough such that  $B_{\delta}(\boldsymbol{x}_*) \subset \mathcal{X}$ ,  $\rho(\delta) < \|D\boldsymbol{f}^{-1}(\boldsymbol{x}_*)\|^{-1}$  and  $D\boldsymbol{f}$  is Lipschitz continuous on  $B_{\delta}(\boldsymbol{x}_*)$  with constant L. Then for  $\boldsymbol{x} \in B_{\delta}(\boldsymbol{x}_*)$  and  $t \in [0, 1]$ ,

$$\|Df(x) - Df(x_* + t(x - x_*))\| \le L \|x - x_* - t(x - x_*)\| = L(1 - t)\|x - x_*\|_{t=0}$$

so the right hand side of (4.1) is not greater than

$$\|D\boldsymbol{f}(\boldsymbol{x})^{-1}\|L\int_{0}^{1}(1-t)\,dt\,\|\boldsymbol{x}-\boldsymbol{x}_{*}\|^{2} = L\|D\boldsymbol{f}(\boldsymbol{x})^{-1}\|\|\boldsymbol{x}-\boldsymbol{x}_{*}\|^{2}/2.$$

Hence for  $x \neq x_*$ ,

$$\frac{\|\boldsymbol{\psi}(\boldsymbol{x}) - \boldsymbol{x}_*\|}{\|\boldsymbol{x} - \boldsymbol{x}_*\|^2} \ \le \ \frac{L\|D\boldsymbol{f}(\boldsymbol{x})^{-1}\|}{2} \ \le \ \frac{L\|D\boldsymbol{f}(\boldsymbol{x}_*)^{-1}\|}{2\big(1 - \|D\boldsymbol{f}(\boldsymbol{x})^{-1}\|\rho(\delta)\big)},$$

and the latter bound converges to  $L \| D \boldsymbol{f}(\boldsymbol{x}_*)^{-1} \| / 2$  as  $\delta \downarrow 0$ .

**Remark 4.5** (Range of attraction). Let  $x_*$  be a zero of f, and let  $\mathcal{X}(x_*)$  be the set of all starting points  $x_0 \in \mathcal{X}$  such that the Newton sequence  $(x_n)_{n\geq 0} = (\psi^n(x_0))_{n\geq 0}$  is well-defined and converges to  $x_*$ . (Here and throughout the sequel,  $\psi^n$  denotes the *n*-fold iteration of  $\psi$ .) The set  $\mathcal{X}(x_*)$  is always an open subset of  $\mathcal{X}$  containing  $x_*$ .

To prove this, note first that by Theorem 4.2,  $\mathcal{X}(\boldsymbol{x}_*)$  contains an open neighborhood  $U_0$  of  $\boldsymbol{x}_*$ . A point  $\boldsymbol{x}_0 \in \mathcal{X} \setminus U_0$  belongs to  $\mathcal{X}(\boldsymbol{x}_*)$  if and only if

for some 
$$n \in \mathbb{N}$$
,  $\boldsymbol{x}_j \in \mathcal{X}$  for  $1 \leq j < n$  and  $\boldsymbol{x}_n \in U_0$ .

Indeed, once the sequence has reached a point in  $U_0$ , it will certainly converge to  $x_*$ . Defining

$$U_n := \psi^{-1}(U_{n-1})$$
 for  $n = 1, 2, 3, \dots$ 

we may reformulate this as  $x_n \in U_n$ . Hence

$$\mathcal{X}(\boldsymbol{x}_*) = \bigcup_{n\geq 0} U_n.$$

But continuity of  $\psi : \mathcal{X} \to \mathbb{R}^d$  and openness of  $U_0$  imply that all sets  $U_0, U_1, U_2, \ldots$  are open. Consequently,  $\mathcal{X}(\boldsymbol{x}_*)$  is an open set, too.

**Example 4.6** (Julia<sup>1</sup> set). Let  $f : \mathbb{C} \to \mathbb{C}$  be given by

$$f(z) := z^3 - 1.$$

90

<sup>&</sup>lt;sup>1</sup>Gaston Maurice Julia, 1893-1978: French mathematician working on rational functions and fractal sets

We may identify a complex number z with the vector  $(\text{Re } z, \text{Im } z)^{\top} \in \mathbb{R}^2$ . Then f is continuously differentiable with (complex) derivative

$$f'(z) = 3z^2$$

corresponding to the Jacobian matrix

$$\begin{bmatrix} \operatorname{Re} f'(z) & -\operatorname{Im} f'(z) \\ \operatorname{Im} f'(z) & \operatorname{Re} f'(z) \end{bmatrix}.$$

Hence we consider  $\mathcal{X} = \mathbb{C} \setminus \{0\}$ . The zeros of f are precisely the third unit roots

$$y_0 := 1, \quad y_1 := \exp(2\pi i/3), \quad y_2 := \exp(4\pi i/3).$$

Newton's method is based on the mapping

$$\psi(z) := z - \frac{f(z)}{f'(z)} = \frac{2z}{3} + \frac{1}{3z^2}$$

According to Remark 4.5, for j = 0, 1, 2 the set

 $\mathcal{X}_j := \{z \in \mathcal{X} : \text{Newton's method with starting point } z \text{ converges to } y_j \}$ 

is an open neighborhood of the point  $y_i$ . One can also verify directly that

$$\{z \in \mathbb{C} : |z - y_j| \le 1/3\} \subset \mathcal{X}_j.$$

The sets  $\mathcal{X}_j$  have a very interesting shape, as shown in Figure 4.1. There  $\mathcal{X}_0$ ,  $\mathcal{X}_1$  and  $\mathcal{X}_2$  are depicted in white, light gray and dark gray, respectively. A remarkable property of these sets is their *fractal* nature: On arbitrarily small scales one can find the same structures as on the original scale. All this is elaborated in the following four exercises.

**Exercise 4.7** (Julia set, I). Let f,  $\mathcal{X}$  and  $\psi$  be as in Example 4.6 with  $f^{-1}(0) = \{y_0, y_1, y_2\}$ , where  $y_j = \exp(2\pi i j/3)$ .

(a) Show that  $f(y_j z) = f(z)$  and  $\psi(y_j z) = y_j \psi(z)$  for arbitrary  $z \in \mathcal{X}$  and j = 0, 1, 2.

(b) Verify that for any  $z \in \mathcal{X}$ , the three equations  $\psi(z) = z$ ,  $\psi'(z) = 0$  and f(z) = 0 are equivalent.

(c) Determine explicit constants C > 0 and  $D \in (0, 1)$  such that for j = 0, 1, 2, 1,

$$|\psi(z) - y_j| \leq D|z - y_j| \quad \text{if } |z - y_j| \leq C.$$

For instance, there exists a solution (C, D) with C = 1/3.

(c') Determine a maximal constant C' > 0 such that for j = 0, 1, 2,

$$|\psi(z) - y_j| < |z - y_j|$$
 if  $|z - y_j| < C'$ .

**Exercise 4.8** (Julia set, II). Now we investigate the inverse image  $\psi^{-1}(y)$  for various  $y \in \mathbb{C}$ . (a) Determine the set  $\psi^{-1}(0) \subset \mathcal{X}$ .



Figure 4.1: Ranges of attraction for Example 4.6.

(**b**) Show that 
$$\psi^{-1}(y_j) = \{y_j, -y_j/2\}$$
 for  $j = 0, 1, 2$ 

(c) Show that for any  $y \in \mathbb{C} \setminus \{y_0, y_1, y_2\}$  the set  $\psi^{-1}(y)$  consists of three different points  $z_1, z_2, z_3 \in \mathcal{X}$ , where

$$\max_{i=1,2,3} |z_i| \ge |y|/2$$
 and  $\min_{i=1,2,3} |z_i| \le |y|^{-1/2}.$ 

**Exercise 4.9** (Julia set, III). From now on we consider the iterates  $\psi^n$ ,  $n \in \mathbb{N}_0$ , of the mapping  $\psi$ . Show that there exist real numbers  $0 = t_0 > t_1 > t_2 > t_3 > \cdots$  with the following properties: (i)  $\psi^n(t_n) = 0$  for  $n \in \mathbb{N}$ .

(ii)  $\lim_{n\to\infty} t_n = -\infty$ .

(iii)  $\lim_{n\to\infty} \psi^n(t) = 1$  for any starting point  $t \in \mathbb{R} \setminus \{0, t_1, t_2, t_3, \ldots\}$ .

Exercise 4.10 (Julia set, IV). Now we combine the results from the previous three exercises.

(a) Deduce from the previous two exercises that 0 is contained in  $\overline{\mathcal{X}(y_0)} \cap \overline{\mathcal{X}(y_1)} \cap \overline{\mathcal{X}(y_2)}$ .

(b) Let  $\psi^{-n}(0)$  be the inverse image  $(\psi^n)^{-1}(0)$ , and let  $\mathcal{X}_o := \bigcup_{n \ge 1} \psi^{-n}(0)$ . That means,  $\mathcal{X}_o$  is the set of problematic starting points for Newton's method.

(b.i) Verify that for arbitrary integers  $1 \le m < n$ ,  $\psi^{-m}(0) \cap \psi^{-n}(0) = \emptyset$ , where  $\psi^{-m}(0)$  consists of  $3^m$  different points in  $\mathcal{X}$ .

(b.ii) Show that  $\mathcal{X}_o = y_j \mathcal{X}_o$  for j = 1, 2.

(b.iii) Show that X<sub>o</sub> contains points z with arbitrarily small and arbitrarily large modulus |z|.
(c) Show that for each integer n ≥ 1 and any point x ∈ ψ<sup>-n</sup>(ℂ),

$$(\psi^n)'(x) = \prod_{k=0}^{n-1} \psi'(\psi^k(x))$$

(with  $\psi^0(x) = x$ ). Deduce from that formula that  $(\psi^n)'(x) \neq 0$  for  $x \in \psi^{-n}(0)$ . Then apply the inverse mapping theorem to show that each point  $x \in \mathcal{X}_o$  belongs to the intersection  $\overline{\mathcal{X}(y_0)} \cap \overline{\mathcal{X}(y_1)} \cap \overline{\mathcal{X}(y_2)}$ . More generally, the intersection of  $\mathcal{X}(y_j)$  with a small neighborhood of x looks approximately like the intersection of  $\mu(x)^{-1}\mathcal{X}(y_j)$  with a small neighborhood of 0.

Exercise 4.11 (Logistic regression, II). As in Example 3.64 let

$$f(\boldsymbol{x}) := \sum_{i=1}^{n} \left( -y_i \, \boldsymbol{z}_i^{\top} \boldsymbol{x} + \log(1 + e^{\boldsymbol{z}_i^{\top} \boldsymbol{x}}) \right)$$

with vectors  $z_1, \ldots, z_n \in \mathbb{R}^d$  such that span $(z_1, \ldots, z_n) = \mathbb{R}^d$  and numbers  $y_1, \ldots, y_n \in \{0, 1\}$ . Write a program for the minimization of f via Newton's method.

**Variations.** Sometimes one replaces the exact derivative Df(x) by an approximation H(x). Two examples for this are:

(A) One approximates the derivative of f by difference ratios. That means, for a small number  $\delta = \delta(x) > 0$ , the Jacobian matrix Df(x) is replaced with

$$oldsymbol{H}(oldsymbol{x}) \ := \ \Big(rac{f_i(oldsymbol{x}+\deltaoldsymbol{e}_j)-f_i(oldsymbol{x}-\deltaoldsymbol{e}_j)}{2\delta}\Big)^d_{i,j=1}.$$

(B) Quite often, the function f considered here is the gradient of a real-valued and twice continuously differentiable function  $\tilde{f}$  to be minimized. In this case, Df(x) is the Hessian matrix of  $\tilde{f}$ at the point x and symmetric. Often the function  $\tilde{f}$  is even strictly convex with Df(x) positive definite for any x. But the condition number of Df(x), the ratio of its largest to its smallest eigenvalue, can be very large, which causes numerical problems when computing  $\psi(x)$ . To avoid this, people often replace Df(x) with

$$oldsymbol{H}(oldsymbol{x}) \ := \ Doldsymbol{f}(oldsymbol{x}) + \epsilon(oldsymbol{x})oldsymbol{I}$$

with some  $\epsilon(\boldsymbol{x}) > 0$ .

Running Newton's method with a surrogate H(x) for Df(x) is called a *quasi-Newton method*. An obvious question is how well it performs. Let us assume that H(x) is always nonsingular and continuous in  $x \in \mathcal{X}$ . Now we consider the algorithmic mapping

$$ilde{oldsymbol{\psi}}(oldsymbol{x}) \ := \ oldsymbol{x} - oldsymbol{H}(oldsymbol{x})^{-1}oldsymbol{f}(oldsymbol{x}).$$

The proof of the following lemma is left to the reader as an exercise:

**Lemma 4.12**. Let  $x_* \in \mathcal{X}$  be a zero of f. Then the algorithmic mapping  $\tilde{\psi}$  defined above has the following property:

$$oldsymbol{\psi}(oldsymbol{x}) - oldsymbol{x}_* \;=\; oldsymbol{A}(oldsymbol{x})(oldsymbol{x} - oldsymbol{x}_*)$$

with a matrix  $oldsymbol{A}(oldsymbol{x}) \in \mathbb{R}^{d imes d}$  such that

$$\lim_{\boldsymbol{x}\to\boldsymbol{x}_*} \boldsymbol{A}(\boldsymbol{x}) = \boldsymbol{I} - \boldsymbol{H}(\boldsymbol{x}_*)^{-1} \boldsymbol{D} \boldsymbol{f}(\boldsymbol{x}_*).$$

This lemma shows that superlinear convergence gets lost in general. Instead,

$$\limsup_{\bm{x} \to \bm{x}_*} \frac{\|\bm{\psi}(\bm{x}) - \bm{x}_*\|}{\|\bm{x} - \bm{x}_*\|} = \|\bm{I} - \bm{H}(\bm{x}_*)^{-1} D \bm{f}(\bm{x}_*)\|$$

If the right hand side is strictly smaller than 1, we still have linear convergence. One can deduce from Lemma 4.3 that

$$\begin{split} \|\boldsymbol{I} - \boldsymbol{H}(\boldsymbol{x}_*)^{-1} D \boldsymbol{f}(\boldsymbol{x}_*)\| &\leq \|\boldsymbol{H}(\boldsymbol{x}_*)^{-1}\| \|\boldsymbol{H}(\boldsymbol{x}_*) - D \boldsymbol{f}(\boldsymbol{x}_*)\| \\ &\leq \frac{\|D \boldsymbol{f}(\boldsymbol{x}_*)^{-1}\| \|\boldsymbol{H}(\boldsymbol{x}_*) - D \boldsymbol{f}(\boldsymbol{x}_*)\|}{1 - \|D \boldsymbol{f}(\boldsymbol{x}_*)^{-1}\| \|\boldsymbol{H}(\boldsymbol{x}_*) - D \boldsymbol{f}(\boldsymbol{x}_*)\|}, \end{split}$$

and the right hand side is strictly smaller than 1 if

$$\| \boldsymbol{H}(\boldsymbol{x}_*) - D\boldsymbol{f}(\boldsymbol{x}_*) \| < \frac{1}{2\| D\boldsymbol{f}(\boldsymbol{x}_*)^{-1} \|}$$

For the special setting (B) with a symmetric and positive definitie matrix  $Df(x_*)$  and  $H(x_*) = Df(x_*) + \epsilon(x_*)I$  one can deduce from the spectral representation of  $Df(x_*)$  that

(4.2) 
$$\|\boldsymbol{I} - \boldsymbol{H}(\boldsymbol{x}_*)^{-1} D\boldsymbol{f}(\boldsymbol{x}_*)\| = \frac{\epsilon(\boldsymbol{x}_*)}{\epsilon(\boldsymbol{x}_*) + \lambda_{\min}(D\boldsymbol{f}(\boldsymbol{x}_*))}$$

with  $\lambda_{\min}(A)$  denoting the smallest real eigenvalue of a matrix A. Thus we may guarantee linear convergence here.

Exercise 4.13. Prove Lemma 4.12.

**Exercise 4.14**. Verify equation (4.2).

## 4.2 Minimization Problems

Now let  $(\mathcal{X}, d)$  be a metric space and  $f : \mathcal{X} \to \mathbb{R}$  be a given function to be minimized. To this end we consider algorithmic mappings  $\psi : \mathcal{X} \to \mathcal{X}$  we intend to iterate. With  $\psi^0(x) := x$  and  $\psi^n(x) := \psi(\psi^{n-1}(x))$  for  $n \in \mathbb{N}$ , we want to find conditions on f and  $\psi$  such that

$$\mathcal{X}_* := \operatorname*{arg\,min}_{\mathcal{X}} f$$

is compact (and thus nonvoid) and any starting point  $x_0 \in \mathcal{X}$  will lead to a sequence  $(x_n)_{n\geq 0} := (\psi^n(x_0))_{n\geq 0}$  converging to the set  $\mathcal{X}_*$ . That means,

$$d(x_n,\mathcal{X}_*) \ := \ \min_{x_*\in\mathcal{X}_*} \ d(x_n,x_*) \ \to \ 0 \quad \text{as} \ n\to\infty.$$

### **4.2.1** A general criterion for convergence

Suppose that f is lower semicontinuous and satisfies the following condition:

(4.3) 
$$\{x \in \mathcal{X} : f(x) \le f(x_o)\} \text{ is compact for arbitrary } x_o \in \mathcal{X}$$

For instance, this condition is satisfied if  $f : \mathbb{R}^d \to \mathbb{R}$  is convex and coercive. Condition (4.3) implies that  $\mathcal{X}_* = \arg \min_{\mathcal{X}}(f)$  is a compact (and thus nonvoid) subset of  $\mathcal{X}$ , see Exercise 3.55.

Now we consider a mapping  $\psi : \mathcal{X} \to \mathcal{X}$  with the following two properties:

(4.4) 
$$f(\psi(x)) = f(x) \quad \text{for } x \in \mathcal{X}_*,$$

(4.5) 
$$f(\psi(x)) < f(x)$$
 for  $x \notin \mathcal{X}_*$ 

These conditions imply that for any starting point  $x_0 \in \mathcal{X}$ , the sequence  $(f(\psi^n(x_0)))_{n\geq 0}$  is nonincreasing. Unfortunately, however, (4.5) is too weak to guarantee convergence of  $(\psi^n(x_0))_{n\geq 0}$  to  $\mathcal{X}_*$ .

**Example 4.15**. Let  $\mathcal{X} = [0, \infty)$  and f(x) = x, so  $\mathcal{X}_* = \{0\}$ . For any  $c \in (0, 1)$ , the algorithmic mapping

$$\psi(x) := \begin{cases} 0 & \text{if } x \le 1 \\ 1 + c(x-1) & \text{if } x > 1 \end{cases}$$

satisfies conditions (4.4) and (4.5). But for any starting point  $x_0 > 1$ ,  $\psi^n(x_0) = 1 + c^n(x_0 - 1) \rightarrow 1$  as  $n \rightarrow \infty$ .

Here is a stronger version of condition (4.5) which is sufficient for our purposes:

(4.6) 
$$\limsup_{x \to y} f(\psi(x)) < f(y) \quad \text{for all } y \in \mathcal{X} \setminus \mathcal{X}_*.$$

For instance, if f and  $\psi$  are continuous, and if  $\psi$  satisfies (4.5), then condition (4.6) is satisfied as well.

**Theorem 4.16**. Suppose that  $\psi$  satisfies conditions (4.4) and (4.6). Then for any starting point  $x_0 \in \mathcal{X}$ , the sequence  $(x_n)_{n>0} = (\psi^n(x_0))_{n>0}$  converges to  $\mathcal{X}_*$ .

**Proof of Theorem 4.16.** As mentioned already, properties (4.4) and (4.5) of  $\psi$  imply that the sequence of numbers  $f(x_n)$  is non-increasing. In particular, all points  $x_n$  are contained in the compact set  $\mathcal{X}_o := \{x \in \mathcal{X} : f(x) \leq f(x_0)\}$ . Suppose that  $(x_n)_n$  does not converge to  $\mathcal{X}_*$ . That means, for some  $\delta > 0$ ,  $d(x_n, \mathcal{X}_*) \geq \delta$  for infinitely many indices n. Since  $\{x \in \mathcal{X}_o : d(x, \mathcal{X}_*) \geq \delta\}$  is compact, too, there exist indices  $n(1) < n(2) < n(3) < \cdots$  such that  $\lim_{k \to \infty} x_{n(k)} = y \in \mathcal{X}_o$  with  $d(y, \mathcal{X}_*) \geq \delta$ . By lower semicontinuity of f and monotonicity of  $(f(x_n))_{n=0}^{\infty}$ ,

$$f(y) \leq \lim_{k \to \infty} f(x_{n(k)})$$
  
= 
$$\lim_{k \to \infty} f(x_{n(k)+1})$$
  
= 
$$\lim_{k \to \infty} f(\psi(x_{n(k)}))$$
  
$$\leq \limsup_{x \to y} f(\psi(x)).$$

Because of condition (4.6), this would imply that  $y \in \mathcal{X}_*$ , a contradiction.

- L	_	_	_	

**Example 4.17** (Coordinatewise descent and LASSO). Let  $\mathcal{X} = \mathbb{R}^d$ , and consider the function  $f_{\lambda}(\boldsymbol{x}) := \|\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{x}\|^2/2 + \lambda \sum_{j=1}^d |x_j|$  with a given vector  $\boldsymbol{y} \in \mathbb{R}^n$ , a matrix  $\boldsymbol{Z} \in \mathbb{R}^{n \times d}$  whose columns have Euclidean norm 1 and some  $\lambda > 0$ . This function is convex and coercive. As shown in Exercise 3.39, for any point  $\boldsymbol{x} \in \mathbb{R}^d$  and index  $k \in \{1, \dots, d\}$  there exists a unique vector

$$\boldsymbol{\psi}^{(k)}(\boldsymbol{x}) := \arg\min\{f_{\lambda}(\tilde{\boldsymbol{x}}) : \tilde{\boldsymbol{x}} \in \mathbb{R}^{d}, \tilde{x}_{j} = x_{j} \text{ whenever } j \neq k\}$$

and this defines a continuous mapping  $\psi^{(k)} : \mathbb{R}^d \to \mathbb{R}^d$ . Moreover, x is a minimizer of  $f_{\lambda}$  if and only if  $\psi^{(k)}(x) = x$  for all  $k \in \{1, \dots, d\}$ . This implies that the algorithmic mapping

$$oldsymbol{\psi} \ := \ oldsymbol{\psi}^{(d)} \circ oldsymbol{\psi}^{(d-1)} \circ \cdots \circ oldsymbol{\psi}^{(1)}$$

is continuous and satisfies conditions (4.4) and (4.5). Since  $f_{\lambda}$  is continuous, too, condition (4.6) is satisfied as well. Consequently, iterating the mapping  $\psi$ , also called iterative coordinatewise descent, yields a sequence converging to  $\arg \min_{x \in \mathbb{R}^d} f_{\lambda}(x)$ .

**Example 4.18** (Spatial median and Weiszfeld's algorithm). As in Exercise 3.37, let  $f : \mathbb{R}^d \to \mathbb{R}$  be given by

$$f(\boldsymbol{x}) := \sum_{i=1}^n \|\boldsymbol{x} - \boldsymbol{x}_i\|$$

with given points  $x_1, \ldots, x_n \in \mathbb{R}^d$ . As shown earlier,

$$Df(\boldsymbol{x}, \boldsymbol{v}) = n(\boldsymbol{x}) \|\boldsymbol{v}\| + \boldsymbol{g}(\boldsymbol{x})^{\top} \boldsymbol{v}$$

with

$$n(m{x}) \ := \ \#\{i:m{x}_i=m{x}\} \quad ext{and} \quad m{g}(m{x}) \ := \ \sum_{i:m{x}_i
eq m{x}} \|m{x}-m{x}_i\|^{-1}(m{x}-m{x}_i).$$

In particular,

$$\min_{\boldsymbol{v}:\|\boldsymbol{v}\|\leq 1} Df(\boldsymbol{x},\boldsymbol{v}) = n(\boldsymbol{x}) - \|\boldsymbol{g}(\boldsymbol{x})\|,$$

so

 $\boldsymbol{x} \in \mathcal{X}_*$  if and only if  $\|\boldsymbol{g}(\boldsymbol{x})\| \leq n(\boldsymbol{x})$ .

An explicit algorithm for the computation of a point in  $\mathcal{X}_*$  has been proposed by Weiszfeld (1937)<sup>2</sup> For a given candidate  $\boldsymbol{x} \in \mathbb{R}^d \setminus \mathcal{X}_*$  we approximate f by a function  $f_{\boldsymbol{x}}$  which is easier to minimize explicitly. Precisely, since  $\|\boldsymbol{a} + \boldsymbol{v}\| \leq \|\boldsymbol{a}\| + (\|\boldsymbol{a} + \boldsymbol{v}\|^2 - \|\boldsymbol{a}\|^2)/(2\|\boldsymbol{a}\|)$  for arbitrary vectors  $\boldsymbol{a} \neq \boldsymbol{0}$  and  $\boldsymbol{v}$  in  $\mathbb{R}^d$ ,

$$f(\boldsymbol{x} + \boldsymbol{v}) = \sum_{i:x_i \neq \boldsymbol{x}} \|\boldsymbol{x} + \boldsymbol{v} - \boldsymbol{x}_i\| + n(\boldsymbol{x}) \|\boldsymbol{v}\|$$
  

$$\leq \sum_{i:x_i \neq \boldsymbol{x}} \|\boldsymbol{x} - \boldsymbol{x}_i\| + \sum_{i:x_i \neq \boldsymbol{x}} \frac{\|\boldsymbol{x} + \boldsymbol{v} - \boldsymbol{x}_i\|^2 - \|\boldsymbol{x} - \boldsymbol{x}_i\|^2}{2\|\boldsymbol{x} - \boldsymbol{x}_i\|} + n(\boldsymbol{x}) \|\boldsymbol{v}\|$$
  

$$= f(\boldsymbol{x}) + \sum_{i:x_i \neq \boldsymbol{x}} \frac{2(\boldsymbol{x} - \boldsymbol{x}_i)^\top \boldsymbol{v} + \|\boldsymbol{v}\|^2}{2\|\boldsymbol{x} - \boldsymbol{x}_i\|} + n(\boldsymbol{x}) \|\boldsymbol{v}\|$$
  

$$= f(\boldsymbol{x}) + \boldsymbol{g}(\boldsymbol{x})^\top \boldsymbol{v} + n(\boldsymbol{x}) \|\boldsymbol{v}\| + w(\boldsymbol{x}) \frac{\|\boldsymbol{v}\|^2}{2}$$
  

$$=: f_{\boldsymbol{x}}(\boldsymbol{x} + \boldsymbol{v})$$

<sup>&</sup>lt;sup>2</sup>Endre Weiszfeld, alias Andrew Vázsonyi (1916-2003): Hungarian mathematician and operations researcher, who emigrated in the 1940's via Paris to the USA.

with  $n(\boldsymbol{x})$  and  $\boldsymbol{g}(\boldsymbol{x})$  as before, and

$$w(oldsymbol{x}) \ := \ \sum_{i \,:\, oldsymbol{x}_i 
eq oldsymbol$$

Note that

$$f_{\boldsymbol{x}} \geq f$$
,  $f_{\boldsymbol{x}}(\boldsymbol{x}) = f(\boldsymbol{x})$  and  $Df_{\boldsymbol{x}}(\boldsymbol{x},\cdot) \equiv Df(\boldsymbol{x},\cdot)$ .

Writing v = ta with a unit vector  $a \in \mathbb{R}^d$  and a scalar t > 0,

$$f_{\boldsymbol{x}}(\boldsymbol{x} + t\boldsymbol{a}) = f(\boldsymbol{x}) + \boldsymbol{g}(\boldsymbol{x})^{\top} \boldsymbol{a} t + n(\boldsymbol{x})t + w(\boldsymbol{x})\frac{t^2}{2} \\ \geq f(\boldsymbol{x}) - (\|\boldsymbol{g}(\boldsymbol{x})\| - n(\boldsymbol{x}))t + w(\boldsymbol{x})\frac{t^2}{2}$$

with equality if and only if  $a = -||g(x)||^{-1}g(x)$ . Then the unique minimizer with respect to t > 0 is given by

$$t = -\frac{\|g(x)\| - n(x)}{w(x)},$$

and the improvement in the value of  $f_x$  equals

$$f_{x}(x) - f_{x}(x + ta) = \frac{(\|g(x)\| - n(x))^{2}}{2w(x)}.$$

Since  $f_x \ge f$  and  $f_x(x) = f(x)$ , the latter quantity is automatically a lower bound for the improvement in the value of f.

For general  $oldsymbol{x} \in \mathbb{R}^d$  we define

$$oldsymbol{\psi}(oldsymbol{x}) \ := \ egin{cases} oldsymbol{x} & ext{if } \|oldsymbol{g}(oldsymbol{x})\| \leq n(oldsymbol{x}), \ oldsymbol{x} - rac{1-n(oldsymbol{x})/\|oldsymbol{g}(oldsymbol{x})\|}{w(oldsymbol{x})} \,oldsymbol{g}(oldsymbol{x}) & ext{if } \|oldsymbol{g}(oldsymbol{x})\| > n(oldsymbol{x}). \end{cases}$$

Then for  $oldsymbol{x} \in \mathbb{R}^d \setminus \mathcal{X}_*,$ 

$$f(m{x}) - f(m{\psi}(m{x})) \ \geq \ rac{ig(\|m{g}(m{x})\| - n(m{x})ig)^2}{2w(m{x})}.$$

Consequently, Weiszfeld's algorithmic mapping  $\psi$  satisfies (4.4) and (4.5).

Unfortunately, (4.6) is violated in general. But one can easily verify that  $\psi$  is continuous on  $\mathbb{R}^d \setminus \{x_1, \ldots, x_n\}$ . Indeed, for  $x \notin \{x_1, \ldots, x_n\}$ ,

$$n(\boldsymbol{x}) = 0,$$
  

$$\boldsymbol{g}(\boldsymbol{x}) = \sum_{i=1}^{n} \|\boldsymbol{x} - \boldsymbol{x}_{i}\|^{-1} (\boldsymbol{x} - \boldsymbol{x}_{i}),$$
  

$$w(\boldsymbol{x}) = \sum_{i=1}^{n} \|\boldsymbol{x} - \boldsymbol{x}_{i}\|^{-1},$$
  

$$\boldsymbol{\psi}(\boldsymbol{x}) = w(\boldsymbol{x})^{-1} \sum_{i=1}^{n} \|\boldsymbol{x} - \boldsymbol{x}_{i}\|^{-1} \boldsymbol{x}_{i},$$
  

$$f(\boldsymbol{x}) - f(\boldsymbol{\psi}(\boldsymbol{x})) \geq \frac{\|\boldsymbol{g}(\boldsymbol{x})\|^{2}}{2w(\boldsymbol{x})}.$$

Thus a possible algorithm with guaranteed convergence runs as follows:

Phase 1: We compute  $f(x_i)$  for  $1 \le i \le n$  and determine a minimizer  $x^0$  of f on  $\{x_1, \ldots, x_n\}$ . Phase 2: If

$$\|\boldsymbol{g}(\boldsymbol{x}^0)\| \leq n(\boldsymbol{x}^0),$$

then  $x^0$  is already a minimizer of f. Otherwise,  $\psi(x^0)$  defines a first point in the open convex set  $\mathcal{X}_0 := \{ x \in \mathbb{R}^d : f(x) < f(x^0) \}$ , which is contained in  $\mathbb{R}^d \setminus \{ x_1, \ldots, x_n \}$ . Since  $\psi$  is continuous on the latter set and satisfies (4.4) and (4.5), the sequence  $(x^k)_{k\geq 0} := (\psi^k(x^0))_{k\geq 0}$ converges automatically to  $\mathcal{X}_*$ .

In the next exercise it will be shown that  $\mathcal{X}_*$  consists of one point  $x_*$ , unless all points  $x_1, \ldots, x_n$  are lying on a straight line in  $\mathbb{R}^d$ .

Exercise 4.19 (Spatial median, II).

(a) For  $x \in \mathbb{R}^d$  let g(x) := ||x||. Show that for  $x \neq 0$ ,

$$abla g(oldsymbol{x}) = oldsymbol{u}$$
 and  $D^2 g(oldsymbol{x}) = \|oldsymbol{x}\|^{-1} (oldsymbol{I}_d - oldsymbol{u}oldsymbol{u}^ op)$ 

with  $u := ||x||^{-1}x$ .

(b) As in Exercise 3.37 let  $f(x) := \sum_{i=1}^{n} ||x - x_i||$  with given points  $x_1, \ldots, x_n \in \mathbb{R}^d$ . Determine the gradient  $\nabla f(x)$  and the Hessian matrix  $D^2 f(x)$  for  $x \in \mathbb{R}^d \setminus \{x_1, \ldots, x_n\}$ . Show that  $D^2 f(x)$  is positive definite, unless all points  $x_1, \ldots, x_n$  are lying on a straight line in  $\mathbb{R}^d$  running through x.

(c) Show that the function f is strictly convex on  $\mathbb{R}^d$ , unless all points  $x_1, \ldots, x_n$  are lying on a straight line in  $\mathbb{R}^d$  (which is always the case if d = 1).

## 4.2.2 Gradient, Newton and quasi-Newton procedures

Now we consider the special case of a continuously differentiable function f on an open subset  $\mathcal{X}$  of  $\mathbb{R}^d$ . We assume condition (4.3) and that

(4.7) 
$$\nabla f(\boldsymbol{x}) \neq 0 \quad \text{for all } \boldsymbol{x} \in \mathcal{X} \setminus \mathcal{X}_*$$

Conditions (4.3) and (4.7) are satisfied if, for instance,  $\mathcal{X}$  is convex and f is convex such that

$$f(\boldsymbol{x}) \to \infty$$
 as  $\boldsymbol{x} \to \partial \mathcal{X}$  or  $\|\boldsymbol{x}\| \to \infty$ .

A model algorithm. In the present context we discuss algorithmic mappings  $\psi : \mathcal{X} \to \mathcal{X}$  of the following type:

$$oldsymbol{\psi}(oldsymbol{x}) \;=\; oldsymbol{x} + \lambda(oldsymbol{x}) \Delta(oldsymbol{x})$$

with a "candidate step function"  $\Delta : \mathcal{X} \to \mathbb{R}^d$  and a "step size function"  $\lambda : \mathcal{X} \to [0,1]$ . The mapping  $\psi$  satisfies the conditions of Theorem 4.16, provided that the following two conditions are fulfilled:

#### (i) $\Delta$ is continuous, and

$$\Delta(\boldsymbol{x}) = \boldsymbol{0} \quad ext{for } \boldsymbol{x} \in \mathcal{X}_*,$$
  
 $abla f(\boldsymbol{x})^\top \Delta(\boldsymbol{x}) < 0 \quad ext{for } \boldsymbol{x} \in \mathcal{X} \setminus \mathcal{X}_*.$ 

(ii) For  $x \in \mathcal{X} \setminus \mathcal{X}_*$  define

$$\begin{split} C(\boldsymbol{x}) &:= \frac{f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x}))}{\max_{t \in [0,1]} \left( f(\boldsymbol{x}) - f(\boldsymbol{x} + t\Delta(\boldsymbol{x})) \right)} \qquad (f := \infty \ \text{on} \ \mathbb{R}^d \setminus \mathcal{X}), \\ \tilde{C}(\boldsymbol{x}) &:= \frac{f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x}))}{-\nabla f(\boldsymbol{x})^\top \Delta(\boldsymbol{x})}. \end{split}$$

Then for any  $\boldsymbol{y} \in \mathcal{X} \setminus \mathcal{X}_*$ ,

$$\liminf_{\boldsymbol{x} \to \boldsymbol{y}} C(\boldsymbol{x}) > 0 \quad \text{or} \quad \liminf_{\boldsymbol{x} \to \boldsymbol{y}} \tilde{C}(\boldsymbol{x}) > 0$$

Concerning the maximum in  $C(\boldsymbol{x})$ , note that  $\{t \in \mathbb{R} : \boldsymbol{x} + t\Delta(\boldsymbol{x}) \in \mathcal{X}\}$  is an open subset of  $\mathbb{R}$  containing 0. Compactness of  $\{f \leq f(\boldsymbol{x})\}$  implies that the set  $\{t \in [0,1] : f(\boldsymbol{x} + t\Delta(\boldsymbol{x})) \leq f(\boldsymbol{x})\}$  is compact, and by condition (4.5),  $\max_{t \in [0,1]} (f(\boldsymbol{x}) - f(\boldsymbol{x} + t\Delta(\boldsymbol{x}))) > 0$ .

To verify (4.6), note that for  $\boldsymbol{y} \in \mathcal{X} \setminus \mathcal{X}_*$ ,

$$\limsup_{\boldsymbol{x} \to \boldsymbol{y}} f(\boldsymbol{\psi}(\boldsymbol{x})) = f(\boldsymbol{y}) - \liminf_{\boldsymbol{x} \to \boldsymbol{y}} (f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x}))),$$

and for a suitable fixed  $t_o = t_o(\boldsymbol{y}) > 0$ ,

$$\begin{aligned} \liminf_{\boldsymbol{x} \to \boldsymbol{y}} \left( f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x})) \right) &\geq \liminf_{\boldsymbol{x} \to \boldsymbol{y}} C(\boldsymbol{x}) \left( f(\boldsymbol{x}) - f(\boldsymbol{x} + t_o\Delta(\boldsymbol{x})) \right) \\ &= \liminf_{\boldsymbol{x} \to \boldsymbol{y}} C(\boldsymbol{x}) \left( f(\boldsymbol{y}) - f(\boldsymbol{y} + t_o\Delta(\boldsymbol{y})) \right) \\ &> 0 \end{aligned}$$

or

$$\begin{split} \liminf_{\boldsymbol{x} \to \boldsymbol{y}} \left( f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x})) \right) &= \liminf_{\boldsymbol{x} \to \boldsymbol{y}} \tilde{C}(\boldsymbol{x}) \left( -\nabla f(\boldsymbol{x})^{\top} \Delta(\boldsymbol{x}) \right) \\ &= \liminf_{\boldsymbol{x} \to \boldsymbol{y}} \tilde{C}(\boldsymbol{x}) \left( -\nabla f(\boldsymbol{y})^{\top} \Delta(\boldsymbol{y}) \right) \\ &> 0. \end{split}$$

## **4.2.3** Examples for the candidate step function $\Delta$

Gradient descent. In the simplest case one chooses

$$\Delta(\boldsymbol{x}) := -\nabla f(\boldsymbol{x}).$$

More generally one could take  $\Delta(\mathbf{x}) := -\kappa(\mathbf{x})\nabla f(\mathbf{x})$  with a continuous mapping  $\kappa : \mathcal{X} \to (0, \infty)$ .

**Newton proposal.** Suppose that f is twice continuously differentiable on  $\mathcal{X}$  with positive definite Hessian matrix  $D^2 f(\boldsymbol{x})$ . Then a promising choice of  $\Delta(\boldsymbol{x})$  seems to be

$$\Delta(\boldsymbol{x}) := -D^2 f(\boldsymbol{x})^{-1} \nabla f(\boldsymbol{x}),$$

see also Section 4.1.

**Quasi-Newton proposal.** Both preceding variants may be generalized and interpreted as follows: For given  $x \in \mathcal{X}$  we approximate f by a quadratic function

$$\boldsymbol{y} \mapsto f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + 2^{-1} (\boldsymbol{y} - \boldsymbol{x})^{\top} \boldsymbol{A}(\boldsymbol{x}) (\boldsymbol{y} - \boldsymbol{x})$$

with a symmetric, positive definite matrix A(x) depending continuously on  $x \in \mathcal{X}$ . The unique minimizer of this quadratic function equals  $x + \Delta(x)$  with

$$\Delta(\boldsymbol{x}) := -\boldsymbol{A}(\boldsymbol{x})^{-1} \nabla f(\boldsymbol{x}).$$

In case of gradient descent, A(x) = I or  $A(x) = \kappa(x)^{-1}I$ , respectively. In case of Newton's proposal,  $A(x) = D^2 f(x)$ . Another possible variant is

$$oldsymbol{A}(oldsymbol{x}) \ := \ \mathrm{diag}\Bigl(\Bigl(rac{\partial^2 f}{\partial x_i^2}(oldsymbol{x})\Bigr)_{i=1}^d\Bigr),$$

provided that the second derivatives on the right hand side exist, are strictly positive and continuous in  $x \in \mathcal{X}$ .

Condition (i) of our model algorithm is always satisfied. Obviously,  $\Delta$  is a continuous function on  $\mathcal{X}$ , and  $\Delta(\boldsymbol{x}) = 0$  if and only if  $\nabla f(\boldsymbol{x}) = 0$  which is equivalent to  $\boldsymbol{x} \in \mathcal{X}_*$ . Moreover,

$$abla f(\boldsymbol{x})^{\top} \Delta(\boldsymbol{x}) = - 
abla f(\boldsymbol{x})^{\top} \boldsymbol{A}(\boldsymbol{x})^{-1} 
abla f(\boldsymbol{x}) < 0$$

for  $x \in \mathcal{X} \setminus \mathcal{X}_*$ .

## **4.2.4** Examples for the step size function $\lambda$

An ideal choice for  $\lambda(x)$  seems to be

$$\lambda(\boldsymbol{x}) \in \operatorname*{arg\,min}_{t \in [0,1]} f(\boldsymbol{x} + t\Delta(\boldsymbol{x})).$$

Obviously condition (ii) in Section 4.2.2 would be satisfied with  $C(\cdot) \equiv 1$ . But in practice one could compute these step sizes  $\lambda(x)$  only approximately, for instance with one of the methods in Chapter 1. Hence one would run two nested iterative algorithms which could be too time-consuming.

With the "improvement function"  $H_{\boldsymbol{x}}: \mathbb{R} \to [-\infty, \infty)$  given by

$$H_{\boldsymbol{x}}(t) := f(\boldsymbol{x}) - f(\boldsymbol{x} + t\Delta(\boldsymbol{x}))$$

$$\lambda(\boldsymbol{x}) \in \underset{t \in [0,1]}{\operatorname{arg\,max}} H_{\boldsymbol{x}}(t).$$

Now we introduce two alternative methods for picking  $\lambda(x)$ . Note that

$$H_{\boldsymbol{x}}(0) = 0 < H'_{\boldsymbol{x}}(0) = -\nabla f(\boldsymbol{x})^{\top} \Delta(\boldsymbol{x}) \text{ for all } \boldsymbol{x} \in \mathcal{X} \setminus \mathcal{X}_{*}.$$

With  $H_x$  one may rewrite C(x) and  $\tilde{C}(x)$  in condition (ii) in Section 4.2.2 as

$$C(\boldsymbol{x}) = rac{H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))}{\sup_{[0,1]}H_{\boldsymbol{x}}} \quad ext{and} \quad ilde{C}(\boldsymbol{x}) = rac{H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))}{H_{\boldsymbol{x}}'(0)}$$

The method of Goldstein–Armijo. We fix a constant  $a \in (0, 1)$  and define

$$\lambda(\mathbf{x}) := 2^{-n(\mathbf{x})}$$
 with  $n(\mathbf{x}) := \min\{n \in \mathbb{N}_0 : H_{\mathbf{x}}(2^{-n}) \ge a 2^{-n} H'_{\mathbf{x}}(0)\}$ 

for  $x \in \mathcal{X} \setminus \mathcal{X}_*$ . Since

$$\lim_{n \to \infty} \frac{H_{\boldsymbol{x}}(2^{-n})}{2^{-n}H'_{\boldsymbol{x}}(0)} = 1 > a$$

the number n(x) is well-defined in  $\mathbb{N}_0$ . To verify Condition (ii) in Section 4.2.2, it suffices to show that for any  $y \in \mathcal{X} \setminus \mathcal{X}_*$ ,

$$\limsup_{\boldsymbol{x} \to \boldsymbol{y}} n(\boldsymbol{x}) < \infty,$$

because

$$\tilde{C}(\boldsymbol{x}) \geq 2^{-n(\boldsymbol{x})} a$$

by definition of  $n(\mathbf{x})$ . Indeed, it follows from  $2^n H_{\mathbf{y}}(2^{-n}) \to H'_{\mathbf{y}}(0)$  that for a suitable  $n_o = n_o(\mathbf{y}) \in \mathbb{N}_0$ ,

$$H_{\mathbf{y}}(2^{-n_o}) > 2^{-n_o} a H'_{\mathbf{y}}(0).$$

Since  $H_{\boldsymbol{x}}(2^{-n_o}) \to H_{\boldsymbol{y}}(2^{-n_o})$  and  $H'_{\boldsymbol{x}}(0) \to H'_{\boldsymbol{y}}(0)$  as  $\boldsymbol{x} \to \boldsymbol{y}$ , this implies that  $n(\boldsymbol{x}) \leq n_o$  if  $\boldsymbol{x}$  is sufficiently close to  $\boldsymbol{y}$ .

The definition of  $\lambda(\boldsymbol{x})$  is illustrated in Figure 4.2. One sees two different examples for a function  $H_{\boldsymbol{x}}$  with  $H_{\boldsymbol{x}}(0) = 0$  and  $H'_{\boldsymbol{x}}(0) = 1$ . In addition the linear functions  $t \mapsto H'_{\boldsymbol{x}}(0)t$ ,  $t \mapsto aH'_{\boldsymbol{x}}(0)t$  with a = 1/3 and the resulting value  $\lambda(\boldsymbol{x})$  are depicted.

An obvious question is what particular value a we should choose. To this end, suppose for the moment that the function  $H = H_x$  is a quadratic and concave function:

$$H(t) = bt - ct^2/2$$

with constants b > 0 and  $c \ge 0$ . Here H'(t) = b - ct, so

$$t_* := \underset{[0,1]}{\operatorname{arg\,max}} H = \underset{[0,1]}{\min(b/c,1)}.$$

Suppose first that  $t_* = 1$ , i.e.  $c \le b$ . Here H is strictly increasing on [0, 1]. Hence it would be desirable to have  $\lambda = \lambda(x) = 1$ , because otherwise  $H(\lambda) \le H(1/2) < H(1)$ . Note that

$$H(1) = b - c/2 \ge b/2$$
 and  $aH'(0) = ab$ .



Figure 4.2: Step size function à la Goldstein-Armijo.

Hence to guarantee  $H(1) \ge aH'(0)$  whenever  $c \le b$ , we should choose

 $a \leq 1/2.$ 

Now suppose that  $t_* < 1$ , i.e. c > b. Then  $t_* = b/c$  and  $H(t_*) = b^2/(2c)$ , so

$$\frac{H(\lambda)}{H(t_*)} = \frac{b\lambda - c\lambda^2/2}{b^2/(2c)} = 2\lambda/t_* - (\lambda/t_*)^2 = 1 - (\lambda/t_* - 1)^2.$$

If  $\lambda < t_*$ , then  $2\lambda \le 1$ , and the inequality  $H(2\lambda) < a2\lambda H'(0)$  reads

$$2b\lambda - 2c\lambda^2 < ab2\lambda.$$

Dividing both sides by  $2b\lambda$  yields the inequality  $1 - \lambda/t_* < a$ , so  $0 < 1 - \lambda/t_* < a$  and

$$\frac{H(\lambda)}{H(t_*)} \ge 1 - a^2.$$

If  $\lambda > t_*$ , then the inequality  $H(\lambda) \ge a\lambda H'(0)$  reads

$$b\lambda - c\lambda^2/2 \ge ab\lambda.$$

Dividing both sides by  $b\lambda$  yields the inequality  $1 - (\lambda/t_*)/2 \ge a$ , so  $1 < \lambda/t_* \le 2(1-a)$  and

$$\frac{H(\lambda)}{H(t_*)} \ge 1 - (1 - 2a)^2.$$

Note that  $1 - a^2$  is decreasing and  $1 - (1 - 2a)^2$  is increasing in  $a \in (0, 1/2]$ . Equating both bounds yields

$$a = 1/3$$
 and  $\frac{H(\lambda)}{H(t_*)} \ge 8/9$ .

More generally, for any value  $a \in [1/4, 1/2]$  we could guarantee that  $H(\lambda)/H(t_*) \ge 3/4$ .

**Coarse binary search for convex functions.** Suppose that f is convex on the convex set  $\mathcal{X}$ , so all functions  $H_x$  are concave. A simple version of the step size function  $\lambda(\cdot)$  is as follows: For  $x \in \mathcal{X} \setminus \mathcal{X}_*$  let

$$\lambda(\boldsymbol{x}) := 2^{-n(\boldsymbol{x})},$$

with n(x) being the smallest number  $n \in \mathbb{N}_0$  satisfying the two inequalities

$$0 < H_{\boldsymbol{x}}(2^{-n}/2) \leq H_{\boldsymbol{x}}(2^{-n}).$$

Figure 4.3 illustrates this definition of  $\lambda(\boldsymbol{x})$ . In both panels one sees  $H_{\boldsymbol{x}}(t)$  for  $t \in [0, 1]$  as well as  $\lambda(\boldsymbol{x})$  and  $\lambda(\boldsymbol{x})/2$ . In the left panel,  $\lambda(\boldsymbol{x}) = 1$ , because  $0 < H_{\boldsymbol{x}}(1/2) \leq H_{\boldsymbol{x}}(1)$ . In the right panel,  $H_{\boldsymbol{x}}(1) < H_{\boldsymbol{x}}(1/2) \geq H_{\boldsymbol{x}}(1/4) > 0$ , whence  $\lambda(\boldsymbol{x}) = 1/2$ .



Figure 4.3: Step size function via coarse binary search.

This step size function fulfills condition (ii) in Section 4.2.2 with  $C(\cdot) \ge 1/2$ : Note first that by concavity of  $H_x$  and construction of  $\lambda(x)$ , the point  $t(x) := \max(\arg \max_{[0,1]} H_x)$  has to fulfill

$$\lambda(\boldsymbol{x})/2 \leq t(\boldsymbol{x}) < 2\lambda(\boldsymbol{x}),$$

because  $0 < H_x(\lambda(x)/2) \le H_x(\lambda(x))$  and  $H_x(\lambda(x)) > H_x(2\lambda(x))$  in case of  $\lambda(x) < 1$ . In case of  $\lambda(x)/2 \le t(x) < \lambda(x)$ , concavity of  $H_x$  implies that

$$\begin{aligned} H_{\boldsymbol{x}}(\lambda(\boldsymbol{x})) &\geq H_{\boldsymbol{x}}(\lambda(\boldsymbol{x})/2) &= H_{\boldsymbol{x}}\left(\left(1 - \frac{\lambda(\boldsymbol{x})/2}{t(\boldsymbol{x})}\right) \cdot 0 + \frac{\lambda(\boldsymbol{x})/2}{t(\boldsymbol{x})} t(\boldsymbol{x})\right) \\ &\geq \left(1 - \frac{\lambda(\boldsymbol{x})/2}{t(\boldsymbol{x})}\right) H_{\boldsymbol{x}}(0) + \frac{\lambda(\boldsymbol{x})/2}{t(\boldsymbol{x})} H_{\boldsymbol{x}}(t(\boldsymbol{x})) \\ &= \frac{\lambda(\boldsymbol{x})}{2t(\boldsymbol{x})} H_{\boldsymbol{x}}(t(\boldsymbol{x})) \\ &\geq \frac{H_{\boldsymbol{x}}(t(\boldsymbol{x}))}{2}, \end{aligned}$$

whence  $H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))/H_{\boldsymbol{x}}(t(\boldsymbol{x})) \geq 1/2.$ 

In case of  $\lambda(\boldsymbol{x}) < t(\boldsymbol{x}) < 2\lambda(\boldsymbol{x})$ , concavity of  $H_{\boldsymbol{x}}$  yields

$$\begin{aligned} H_{\boldsymbol{x}}(\lambda(\boldsymbol{x})) &= H_{\boldsymbol{x}}\left(\left(1 - \frac{\lambda(\boldsymbol{x})}{t(\boldsymbol{x})}\right) \cdot 0 + \frac{\lambda(\boldsymbol{x})}{t(\boldsymbol{x})}t(\boldsymbol{x})\right) \\ &\geq \left(1 - \frac{\lambda(\boldsymbol{x})}{t(\boldsymbol{x})}\right)H_{\boldsymbol{x}}(0) + \frac{\lambda(\boldsymbol{x})}{t(\boldsymbol{x})}H_{\boldsymbol{x}}(t(\boldsymbol{x})) \\ &= \frac{\lambda(\boldsymbol{x})}{t(\boldsymbol{x})}H_{\boldsymbol{x}}(t(\boldsymbol{x})) \\ &> \frac{H_{\boldsymbol{x}}(t(\boldsymbol{x}))}{2}, \end{aligned}$$

so again  $H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))/H_{\boldsymbol{x}}(t(\boldsymbol{x})) \geq 1/2$ .

It is also interesting to consider once more the special case of a quadratic and concave function  $H = H_x$ , i.e.  $H(t) = bt - ct^2/2$  with b > 0 and  $c \ge 0$ . Let  $\lambda = \lambda(x)$  and  $t_* := \min\{b/c, 1\}$ , the unique maximizer of H on [0, 1]. In case of  $c \le b$ , the function H is strictly increasing on [0, 1], so  $\lambda = t_*$ . In case of c > b,  $t_* \in (0, 1)$ , and

$$\frac{H(\lambda)}{H(t_*)} = 1 - \left(\frac{\lambda}{t_*} - 1\right)^2,$$

see the considerations for the Goldstein–Armijo method. But the symmetry of H around  $t_*$  together with the inequalities  $H(\lambda/2) \leq H(\lambda)$  and  $H(2\lambda) < H(\lambda)$  in case of  $\lambda < 1$  imply that even

$$\frac{3}{4}\lambda \leq t_* \leq \frac{3}{2}\lambda.$$

Hence  $2/3 \leq \lambda/t_* \leq 4/3$  and

$$\frac{H(\lambda)}{H(t_*)} = 1 - \left(\frac{\lambda}{t_*} - 1\right)^2 \ge \frac{8}{9}.$$

Implementation of the quasi-Newton method with stepsize correction. The actual implementation of a quasi-Newton method with one of the former two stepsize corrections is rather simple, and one does not need to work with the auxiliary index n(x) explicitly. Here is pseudo-code for a quasi-Newton method with the Goldstein–Armijo stepsize correction, where we assume that we have access to the functions  $f : \mathcal{X} \to \mathbb{R}, \nabla f : \mathcal{X} \to \mathbb{R}^d$  and  $A : \mathcal{X} \to \mathbb{R}^{d \times d}_{sym,+}$ , and  $\delta_o > 0$  is a small threshold for the stopping criterion:

$$\begin{array}{l} \boldsymbol{x} \leftarrow \boldsymbol{x}_{o} \\ \boldsymbol{f} \leftarrow \boldsymbol{f}(\boldsymbol{x}) \\ \boldsymbol{g} \leftarrow \nabla \boldsymbol{f}(\boldsymbol{x}) \\ \Delta \leftarrow \boldsymbol{A}(\boldsymbol{x})^{-1} \boldsymbol{g} \\ \delta \leftarrow \boldsymbol{g}^{\top} \Delta \\ \textbf{while } \delta > \delta_{o} \ \textbf{do} \\ \boldsymbol{x}_{new} \leftarrow \boldsymbol{x} - \Delta \\ \boldsymbol{f}_{new} \leftarrow \boldsymbol{f}(\boldsymbol{x}_{new}) \\ \textbf{while } \boldsymbol{f} - \boldsymbol{f}_{new} < \delta/3 \ \textbf{do} \\ \boldsymbol{x}_{new} \leftarrow (\boldsymbol{x} + \boldsymbol{x}_{new})/2 \\ \delta \leftarrow \delta/2 \\ \boldsymbol{f}_{new} \leftarrow \boldsymbol{f}(\boldsymbol{x}_{new}) \\ \textbf{end while} \\ \boldsymbol{x} \leftarrow \boldsymbol{x}_{new} \\ \boldsymbol{f} \leftarrow \boldsymbol{f}_{new} \\ \boldsymbol{g} \leftarrow \nabla \boldsymbol{f}(\boldsymbol{x}) \\ \Delta \leftarrow \boldsymbol{A}(\boldsymbol{x})^{-1} \boldsymbol{g} \\ \delta \leftarrow \boldsymbol{g}^{\top} \Delta \\ \textbf{end while.} \end{array}$$

And here is pseudo-code for a variant with coarse binary search:

$$\begin{array}{l} \boldsymbol{x} \leftarrow \boldsymbol{x}_{o} \\ \boldsymbol{f} \leftarrow \boldsymbol{f}(\boldsymbol{x}) \\ \boldsymbol{g} \leftarrow \nabla \boldsymbol{f}(\boldsymbol{x}) \\ \Delta \leftarrow \boldsymbol{A}(\boldsymbol{x})^{-1} \boldsymbol{g} \\ \delta \leftarrow \boldsymbol{g}^{\top} \Delta \\ \textbf{while } \delta > \delta_{o} \ \textbf{do} \\ \boldsymbol{x}_{new} \leftarrow \boldsymbol{x} - \Delta \\ \boldsymbol{f}_{new} \leftarrow \boldsymbol{f}(\boldsymbol{x}_{new}) \\ \boldsymbol{x}'_{new} \leftarrow (\boldsymbol{x} + \boldsymbol{x}_{new})/2 \\ \boldsymbol{f}'_{new} \leftarrow \boldsymbol{f}(\boldsymbol{x}'_{new}) \\ \textbf{while } \boldsymbol{f}_{new} \geq \boldsymbol{f} \ \textbf{or} \ \boldsymbol{f}'_{new} < \boldsymbol{f}_{new} \ \textbf{do} \\ \boldsymbol{x}_{new} \leftarrow \boldsymbol{x}'_{new} \\ \boldsymbol{f}_{new} \leftarrow \boldsymbol{f}'_{new} \\ \boldsymbol{x}'_{new} \leftarrow (\boldsymbol{x} + \boldsymbol{x}_{new})/2 \\ \boldsymbol{f}'_{new} \leftarrow \boldsymbol{f}(\boldsymbol{x}'_{new}) \\ \textbf{end while} \\ \boldsymbol{x} \leftarrow \boldsymbol{x}_{new} \\ \boldsymbol{f} \leftarrow \boldsymbol{f}(\boldsymbol{x}'_{new}) \\ \textbf{end while} \\ \boldsymbol{g} \leftarrow \nabla \boldsymbol{f}(\boldsymbol{x}) \\ \Delta \leftarrow \boldsymbol{A}(\boldsymbol{x})^{-1} \boldsymbol{g} \\ \delta \leftarrow \boldsymbol{g}^{\top} \Delta \\ \textbf{end while.} \end{array}$$

In practice, one should secure the inner while loop with a counter to make sure that it is not repeated endlessly due to numerical errors. And right after the inner while loop, one should include an extra check whether really  $f_{\text{new}} < f$ . If not, one should just set  $\delta \leftarrow 0$  and thus stop the outer while loop.

**Exercise 4.20** (Another step size function, I). Suppose that f is convex and continuously differentiable on the convex open set  $\mathcal{X}$ , and let  $\Delta : \mathcal{X} \to \mathbb{R}^d$  be a candidate step function with the stated properties. For  $x \in \mathcal{X} \setminus \mathcal{X}_*$  define

$$\lambda(\boldsymbol{x}) := 2^{-n(\boldsymbol{x})}$$

with

$$n(\mathbf{x}) := \min\{n \in \mathbb{N}_0 : H_{\mathbf{x}}(2 \cdot 2^{-n}) \ge 0\}.$$

(a) Show that this stepfunction satisfies the inequality

$$\frac{H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))}{\max_{[0,1]}H_{\boldsymbol{x}}} \geq \frac{1}{4}.$$

(b) Show that in case of a quadratic function  $H_x$ ,

$$\frac{H_{\boldsymbol{x}}(\lambda(\boldsymbol{x}))}{\max_{[0,1]}H_{\boldsymbol{x}}} \geq \frac{3}{4}.$$

#### 4.2.5 Performance in connection with quasi-Newton methods

If we determine  $\Delta(x)$  via a quasi-Newton method with good local convergence properties, it is desirable to have  $\lambda(x) = 1$  for  $x \in \mathcal{X} \setminus \mathcal{X}_*$  sufficiently close to  $\mathcal{X}_*$ . Suppose that f is twice continuously differentiable and convex. Then  $H_x$  is real-valued and twice differentiable on [0, 1]whenever x is sufficiently close to  $\mathcal{X}_*$ . In that case, both the Goldstein–Armijo method with  $a \leq 1/3$  and the rough binary search yield

$$\lambda(\boldsymbol{x}) = 1$$
 whenever  $H''_{\boldsymbol{x}} \ge -\frac{4}{3}H'_{\boldsymbol{x}}(0)$  on  $[0,1]$ .

To verify this claim, note that

$$H_{\boldsymbol{x}}(1) = H'_{\boldsymbol{x}}(0) + \int_{0}^{1} (H'_{\boldsymbol{x}}(t) - H'_{\boldsymbol{x}}(0)) dt$$
$$\geq H'_{\boldsymbol{x}}(0) - \frac{4}{3} H'_{\boldsymbol{x}}(0) \int_{0}^{1} t dt$$
$$= \frac{H'_{\boldsymbol{x}}(0)}{3}$$

and

$$\begin{aligned} H_{\boldsymbol{x}}(1) - H_{\boldsymbol{x}}(1/2) &= \frac{H_{\boldsymbol{x}}'(0)}{2} + \int_{1/2}^{1} (H_{\boldsymbol{x}}'(t) - H_{\boldsymbol{x}}'(0)) \, dt \\ &\geq \frac{H_{\boldsymbol{x}}'(0)}{2} - \frac{4}{3} \, H_{\boldsymbol{x}}'(0) \int_{1/2}^{1} t \, dt \\ &= 0, \end{aligned}$$

because  $H'_{x}(t) - H'_{x}(0) = H''_{x}(\xi(t))t \ge -(4/3)H'_{x}(0)t$  for some point  $\xi(t) \in [0, t]$ .

Specifically let  $\Delta(x) = -A(x)^{-1} \nabla f(x)$  with a symmetric and positive definite matrix A(x) depending continuously on  $x \in \mathcal{X}$ . For  $x_* \in \mathcal{X}_*$  and  $x \in \mathcal{X} \setminus \mathcal{X}_*$  sufficiently close to  $x_*$ , the whole line segment  $\{x + t\Delta(x) : t \in [0, 1]\}$  is contained in  $\mathcal{X}$ , and

$$H'_{\boldsymbol{x}}(0) = -\nabla f(\boldsymbol{x})^{\top} \Delta(\boldsymbol{x})$$
  
=  $\Delta(\boldsymbol{x})^{\top} A(\boldsymbol{x}) \Delta(\boldsymbol{x}),$   
 $H''_{\boldsymbol{x}}(t) = -\Delta(\boldsymbol{x})^{\top} D^2 f(\boldsymbol{x} + t \Delta(\boldsymbol{x})) \Delta(\boldsymbol{x})$ 

Hence

$$H_{\boldsymbol{x}}''(t) + \frac{4}{3}H_{\boldsymbol{x}}'(0) = \Delta(\boldsymbol{x})^{\top} \Big(\frac{4}{3}\boldsymbol{A}(\boldsymbol{x}) - D^2 f(\boldsymbol{x} + t\Delta(\boldsymbol{x}))\Big) \Delta(\boldsymbol{x}) \geq 0,$$

provided that  $(4/3)\mathbf{A}(\mathbf{x}) - D^2 f(\mathbf{x} + t\Delta(\mathbf{x}))$  is positive semidefinite. Consequently,  $\lambda(\mathbf{x}) = 1$  for  $\mathbf{x} \in \mathcal{X} \setminus \mathcal{X}_*$  sufficiently close to  $\mathbf{x}_*$ , provided that

$$\boldsymbol{A}(\boldsymbol{x}_*) - \frac{3}{4} D^2 f(\boldsymbol{x}_*)$$
 is positive definite.

In particular, let f be twice continuously differentiable on the convex open set  $\mathcal{X}$  with positive definite Hessian matrix  $D^2 f(\boldsymbol{x})$  for all  $\boldsymbol{x} \in \mathcal{X}$ , and let  $\Delta(\boldsymbol{x})$  be the Newton step  $-D^2 f(\boldsymbol{x})^{-1} \nabla f(\boldsymbol{x})$ . Then iterating the algorithmic mapping  $\boldsymbol{\psi}$  with arbitrary starting point  $\boldsymbol{x}_0$  will yield a sequence  $(\boldsymbol{x}_n)_{n\geq 0}$  with limit  $\mathcal{X}_* = \{\boldsymbol{x}_*\}$ , and  $\lambda(\boldsymbol{x}_n) < 1$  for at most finitely many n. Hence the step size correction yields guaranteed convergence without sacrificing the rapid convergence of Newton's method.

**Exercise 4.21** (Another step size function, II). Consider the step size function  $\lambda(\cdot)$  introduced in Exercise 4.20. Suppose that  $H_x$  is twice differentiable on [0, 2] with  $H''_x \ge -CH'_x(0)$  for some constant C > 0. Which value of C guarantees that  $\lambda(x) = 1$ ? Would you recommend this step size function in connection with the usual Newton candidate step function?

**Exercise 4.22** (Logistic regression, III). Complement your Newton procedure for logistic regression in Exercise 4.11 with a step size function of your choice.

**Exercise 4.23** (Smooth approximations). Several of our examples involved convex but non-differentiable functions  $\boldsymbol{x} \mapsto \|\boldsymbol{x} - \boldsymbol{x}_o\|$  or  $\boldsymbol{x} \mapsto |\boldsymbol{x}_o^\top \boldsymbol{x} - r_o|$ . A simple way to avoid the problems resulting from this lack of smoothness is to approximate these functions by smooth ones.

(a) For  $\epsilon > 0$  let  $h_{\epsilon} : \mathbb{R} \to \mathbb{R}$  be given by

$$h_{\epsilon}(x) := \sqrt{x^2 + \epsilon^2}.$$

Show that

$$|x| \leq h_{\epsilon}(x) \leq |x| + \min\left(\epsilon, \frac{\epsilon^2}{2|x|}\right),$$

and

$$egin{array}{rl} h'_{\epsilon}(r) &=& rac{r}{\sqrt{\epsilon^2+r^2}}, \ h''_{\epsilon}(r) &=& rac{\epsilon^2}{(\epsilon^2+r^2)^{3/2}}. \end{array}$$

(b) Now let  $h_{\epsilon}(\boldsymbol{x}) := \sqrt{\|\boldsymbol{x}\|^2 + \epsilon^2}$ . Determine the gradient and Hessian matrix of  $h_{\epsilon}$ . Show that the Hessian matrix is always positive definite.

(c) Design and implement a Newton method with regularization of the Hessian matrix and some step size correction for the (approximate) minimization of

$$f(oldsymbol{x}) \ := \ \sum_{i=1}^n |oldsymbol{z}_i^ op oldsymbol{x} - y_i|$$

with given numbers  $y_1, \ldots, y_n \in \mathbb{R}$  and vectors  $z_1, \ldots, z_n \in \mathbb{R}^d$ . Choose a parameter  $\epsilon > 0$  which is related to the numbers  $y_1, \ldots, y_n$  in a reasonable fashion.

**Exercise 4.24**. Let  $X_1, X_2, \ldots, X_n$  be stochastically independent random variables with values in  $\mathbb{R}$  and density function  $f_{\mu_*,\sigma_*}$ . Here  $\mu_* \in \mathbb{R}$  and  $\sigma_* > 0$  are unknown parameters, and

$$f_{\mu,\sigma}(x) = \frac{1}{\sigma} f_{0,1}\left(\frac{x-\mu}{\sigma}\right)$$
 with  $f_{0,1}(y) := \frac{1}{\pi(1+y^2)}$ 

Now we'd like to compute a maximum-likelihood estimator for the unknown parameter vector  $\boldsymbol{\theta}_* := (\mu_*, \sigma_*)^\top$ . That is a parameter vector  $\boldsymbol{\theta} = (\mu, \sigma)^\top \in \mathbb{R} \times (0, \infty)$  minimizing the negative log-likelihood

$$\widehat{L}(\boldsymbol{\theta}) := -\sum_{i=1}^n \log f_{\mu,\sigma}(X_i).$$

Implement a Newton method with step size correction for this task.
# Chapter 5

# **Constrained Optimization**

# 5.1 Lagrange Multipliers

# 5.1.1 The general principle

We consider functions  $f : \mathcal{X} \to (-\infty, \infty]$  and  $g : \mathcal{X} \to \mathbb{R}^q$  on an arbitrary set  $\mathcal{X}$ . For a given vector  $c \in \mathbb{R}^q$  we would like to minimize f under the constraint that g = c or that  $g \leq c$  component-wise.

In many cases it is much easier to minimize the function  $f - \lambda^{\top} g$  on  $\mathcal{X}$  for an arbitrary vector  $\lambda \in \mathbb{R}^q$  and then to study the minimizer as a function of  $\lambda$ . The next theorem explains to what extent this approach works.

**Theorem 5.1**. For some  $\lambda \in \mathbb{R}^q$  let

$$x_{\lambda} \in \operatorname*{arg\,min}_{x \in \mathcal{X}} \left( f(x) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(x) \right)$$

Then

$$x_{\lambda} \in \underset{x \in \mathcal{X}: g(x) = g(x_{\lambda})}{\operatorname{arg\,min}} f(x).$$

More generally, let

$$\mathcal{X}_{\lambda} := \left\{ x \in \mathcal{X} : \boldsymbol{\lambda}^{\top} \boldsymbol{g}(x) \geq \boldsymbol{\lambda}^{\top} \boldsymbol{g}(x_{\lambda}) \right\}$$

Then

$$x_{\lambda} \in \operatorname*{arg\,min}_{x \in \tilde{\mathcal{X}}} f(x)$$

for any set  $\tilde{\mathcal{X}}$  such that  $x_{\lambda} \in \tilde{\mathcal{X}} \subset \mathcal{X}_{\lambda}$ .

If  $x_{\lambda}$  is the unique minimizer of  $f - \lambda^{\top} g$  on  $\mathcal{X}$ , then  $x_{\lambda}$  is also the unique minimizer of f on any set  $\tilde{\mathcal{X}}$  such that  $x_{\lambda} \in \tilde{\mathcal{X}} \subset \mathcal{X}_{\lambda}$ .

**Remark 5.2** (Equality constraints). For the original problem, minimizing f under the constraint that g = c, we obtain the following procedure: For  $\lambda \in \mathbb{R}^q$  we determine a minimizer  $x_{\lambda}$  of  $f - \lambda^{\top} g$  over  $\mathcal{X}$ . Then we try to find  $\lambda \in \mathbb{R}^q$  such that  $g(x_{\lambda}) = c$ . Then the corresponding point  $x_{\lambda}$  solves the original problem.

110

**Remark 5.3** (Inequality constraints). Suppose for a vector  $c = (c_j)_{j=1}^q \in \mathbb{R}^q$  and some index  $q_o \in \{1, \ldots, q\}$  we want to minimize f(x) under the following constraints:

(5.1) 
$$g_j(x) \begin{cases} \leq c_j & \text{for } 1 \leq j \leq q_o, \\ = c_j & \text{for } q_o < j \leq q. \end{cases}$$

For this purpose we try to find a vector  $\lambda \in \mathbb{R}^q$  such that the corresponding minimizer  $x_\lambda$  of  $f - \lambda^\top g$  satisfies (5.1), and

$$\lambda_j \begin{cases} \leq 0 & \text{if } j \leq q_o, \\ = 0 & \text{if } j \leq q_o \text{ and } g_j(x_\lambda) < c_\lambda. \end{cases}$$

If such a vector  $\lambda$  exists, the corresponding  $x_{\lambda}$  solves the given constrained optimization problem, because any point  $x \in \mathcal{X}$  satisfying (5.1) satisfies also  $\lambda_j g_j(x) \ge \lambda_j g_j(x_{\lambda})$  for  $1 \le j \le q$ , so  $x \in \mathcal{X}_{\lambda}$ .

**Proof of Theorem 5.1.** Let y be an arbitrary point in  $\mathcal{X}_{\lambda}$  such that  $f(y) \leq f(x_{\lambda})$ . Then

$$f(x_{\lambda}) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(x_{\lambda}) \geq f(y) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(x_{\lambda}) \geq f(y) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(y),$$

where the last inequality follows from the definition of  $\mathcal{X}_{\lambda}$ . Optimality of  $x_{\lambda}$  implies that both inequalities are equalities, so  $f(y) = f(x_{\lambda})$  and  $\lambda^{\top} g(y) = \lambda^{\top} g(x_{\lambda})$ . If  $x_{\lambda}$  is even the unique minimizer of  $f - \lambda^{\top} g$  on  $\mathcal{X}$ , then the inequality  $f(y) \leq f(x_{\lambda})$  implies that  $y = x_{\lambda}$ .

**Exercise 5.4** (Varying Lagrange multipliers). For two different vectors  $\lambda_1, \lambda_2 \in \mathbb{R}^q$  let

$$x_i \in \underset{x \in \mathcal{X}}{\operatorname{arg\,min}}(f - \boldsymbol{\lambda}_i^{\top} \boldsymbol{g})(x), \quad i = 1, 2.$$

Show that

$$(\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^{\top} (\boldsymbol{g}(x_1) - \boldsymbol{g}(x_2)) \geq 0.$$

Show also that this inequality is strict whenever the two minimizers  $x_1$  and  $x_2$  are unique and different.

# 5.1.2 Examples for Lagrange's Method

**Optimal linear combinations of estimators.** Let  $X_1, X_2, \ldots, X_d$  be stochastically independent random variables with one and the same unknown expected value  $\mu \in \mathbb{R}$  and known standard deviations  $\sigma_i = \text{Std}(X_i) > 0$ . (An example are *d* different measurement devices with known imprecision, and with each of them one measures a sample from one and the same substance.) Now we want to estimate  $\mu$  by a weighted average of the  $X_i$ ,

$$\widehat{\mu} := \sum_{i=1}^d w_i X_i$$

with certain weights  $w_i \in \mathbb{R}$ . A natural constraint on  $\boldsymbol{w} = (w_i)_{i=1}^d$  is that

$$g(\boldsymbol{w}) := \sum_{i=1}^d w_i = 1,$$

because then  $\mathbb{E}(\hat{\mu}) = \mu$ . Now we would like to minimize the mean quadratic error

$$\mathbb{E}\left((\widehat{\mu}-\mu)^2\right) = f(\boldsymbol{w}) := \sum_{i=1}^d w_i^2 \sigma_i^2$$

under that constraint. To this end we minimize for an arbitrary number  $\lambda \in \mathbb{R}$  the function

$$f(\boldsymbol{w}) - \lambda g(\boldsymbol{w}) = \sum_{i=1}^{d} (w_i^2 \sigma_i^2 - \lambda w_i).$$

This minimization may be performed component-wise, and the unique minimizer is given by

$$oldsymbol{w}_{\lambda} \ := \ \Big(rac{\lambda}{2\sigma_i^2}\Big)_{i=1}^d.$$

The constraint g(w) = 1 is satisfied if and only if  $\lambda = 2C$  with

$$C := \left(\sum_{i=1}^{d} \frac{1}{\sigma_i^2}\right)^{-1}.$$

Hence optimal weights are given by

$$w_i := C/\sigma_i^2,$$

and the resulting mean squared error equals C.

Minimizing a quadratic function under linear constraints. The previous example may be considered as a special case of the following optimization problem: Let  $A \in \mathbb{R}^{d \times d}$  be symmetric and positive definite. Now we would like to minimize

$$f(\boldsymbol{x}) := 2^{-1} \boldsymbol{x}^{\top} \boldsymbol{A} \boldsymbol{x}$$

over all  $\boldsymbol{x} \in \mathbb{R}^d$  satisfying the constraint

$$B^{\top}x = c$$

Here **B** is a given matrix in  $\mathbb{R}^{d \times q}$  of rank q < d, and **c** is a given vector in  $\mathbb{R}^{q}$ .

To this end we minimize for  $\lambda \in \mathbb{R}^q$  the function  $f(x) - \lambda^\top B^\top x = f(x) - (B\lambda)^\top x$ . Note that the gradient and Hessian matrix of the latter function are given by  $Ax - B\lambda$  and A, respectively. Thus its unique minimizer is given by

$$x_\lambda \ := \ A^{-1}B\lambda$$
 .

Alternatively one could argue with quadratic completion:  $f(x) - (B\lambda)^{\top} x$  equals

$$2^{-1}\boldsymbol{x}^{\top}\boldsymbol{A}\boldsymbol{x} - (\boldsymbol{B}\boldsymbol{\lambda})^{\top}\boldsymbol{x} = 2^{-1} \Big(\boldsymbol{x}^{\top}\boldsymbol{A}\boldsymbol{x} - 2(\boldsymbol{A}^{-1}\boldsymbol{B}\boldsymbol{\lambda})^{\top}\boldsymbol{A}\boldsymbol{x}\Big)$$
  
=  $2^{-1}(\boldsymbol{x} - \boldsymbol{A}^{-1}\boldsymbol{B}\boldsymbol{\lambda})^{\top}\boldsymbol{A}(\boldsymbol{x} - \boldsymbol{A}^{-1}\boldsymbol{B}\boldsymbol{\lambda}) - 2^{-1}\boldsymbol{\lambda}^{\top}\boldsymbol{B}^{\top}\boldsymbol{A}^{-1}\boldsymbol{B}\boldsymbol{\lambda}.$ 

Furthermore,

$$\boldsymbol{B}^{\top}\boldsymbol{x}_{\lambda} = \boldsymbol{B}^{\top}\boldsymbol{A}^{-1}\boldsymbol{B}\boldsymbol{\lambda},$$

and this equals c if and only if

$$\boldsymbol{\lambda} = (\boldsymbol{B}^{\top} \boldsymbol{A}^{-1} \boldsymbol{B})^{-1} \boldsymbol{c}.$$

Consequently the original constrained optimization problem has the unique solution

$$m{x}_* \ := \ m{A}^{-1} m{B} (m{B}^ op m{A}^{-1} m{B})^{-1} m{c}$$

with

$$f(\boldsymbol{x}_*) = 2^{-1} \boldsymbol{c}^\top (\boldsymbol{B}^\top \boldsymbol{A}^{-1} \boldsymbol{B})^{-1} \boldsymbol{c}.$$

**Exercise 5.5** (Minimizing a quadratic function, I). Generalize the previous considerations to the case of  $f(\boldsymbol{x}) = \gamma + \boldsymbol{\beta}^{\top} \boldsymbol{x} + 2^{-1} \boldsymbol{x}^{\top} \boldsymbol{A} \boldsymbol{x}$  with  $\gamma \in \mathbb{R}$  and  $\boldsymbol{\beta} \in \mathbb{R}^d$ .

Exercise 5.6 (Minimizing a quadratic function, II). (a) Minimize

$$f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^{\top} \boldsymbol{A} \boldsymbol{x} + \boldsymbol{\beta}^{\top} \boldsymbol{x}$$

over all  $oldsymbol{x} \in \mathbb{R}^3$  satisfying

$$\mathbf{b}^{\mathsf{T}} \mathbf{x} = 4,$$

where

$$\boldsymbol{A} = \begin{bmatrix} 3 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad \boldsymbol{b} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}.$$

(b) What happens to the solution in part (a) if the constraint (5.2) is replaced with

$$\boldsymbol{b}^{\top} \boldsymbol{x} \leq (\geq) 4?$$

**Density functions with maximal entropy.** For a probability density p with respect to Lebesgue measure on a Borel set  $\mathcal{X} \subset \mathbb{R}^d$  with  $\text{Leb}(\mathcal{X}) > 0$ , the (differential Shannon) entropy

$$H(p) := -\int_{\mathcal{X}} p(\boldsymbol{x}) \log p(\boldsymbol{x}) \, d\boldsymbol{x}$$

measures how "diffuse" the corresponding distribution is; higher values of H(p) indicate a more diffuse distribution. An obvious question is what the most diffuse distribution looks like. Note that H(p) is well-defined in  $[-\infty, \infty)$  for arbitrary integrable functions  $p : \mathcal{X} \to [0, \infty)$ , with the convention that  $0 \log 0 := 0$ , because  $t \mapsto t \log t$  is bounded from below on  $[0, \infty)$ .

More formally, we consider the set  $\mathcal{P}$  of all nonnegative functions  $p \in L^1(\mathcal{X})$  and the functional  $f(p) := -H(p) \in (0, \infty]$ . The additional constraint that  $p \in \mathcal{P}$  is a probability density corresponds to the requirement

$$g_0(p) := \int_{\mathcal{X}} p(\boldsymbol{x}) \, d\boldsymbol{x} \stackrel{!}{=} 1.$$

Case 1: No further constraints, but  $\text{Leb}(\mathcal{X}) < \infty$ . We want to maximize H(p) = -f(p) among all  $p \in \mathcal{X}$  with  $g_0(p) = 1$ . To this end we minimize

$$f(p) - \lambda g_0(p) = \int_{\mathcal{X}} \left( p(\boldsymbol{x}) \log p(\boldsymbol{x}) - \lambda p(\boldsymbol{x}) \right) d\boldsymbol{x}.$$

Since

$$\frac{d}{dt}\left(t\log t - \lambda t\right) = \log t + 1 - \lambda,$$

the minimizing density is unique (almost everywhere) and given by

$$p_{\lambda} \equiv \exp(\lambda - 1).$$

With  $\lambda := 1 - \log \operatorname{Leb}(\mathcal{X})$  we obtain the probability density  $p_* \equiv \operatorname{Leb}(\mathcal{X})^{-1}$ , the density of the uniform distribution on  $\mathcal{X}$  with  $H(p) = \log \operatorname{Leb}(\mathcal{X})$ .

Case 2: Conditions on the first and second moments. In case of  $\operatorname{Leb}(\mathcal{X}) = \infty$ , there exists no maximizer  $p \in \mathcal{P}$  of H under the constraint  $g_0(p) = 1$ . This can be verified by considering the uniform distribution on  $\mathcal{X}^{(R)} := \{ \boldsymbol{x} \in \mathcal{X} : \|\boldsymbol{x}\| \leq R \}$  for sufficiently large radius R > 0 with density  $p^{(R)}(\boldsymbol{x}) := \operatorname{Leb}(\mathcal{X}^{(R)})^{-1} \mathbf{1}_{\mathcal{X}^{(R)}}(\boldsymbol{x})$  and entropy  $H(p^{(R)}) = \log \operatorname{Leb}(\mathcal{X}^{(R)}) \to \infty$  as  $R \to \infty$ .

Suppose that  $\mathcal{X} = \mathbb{R}^d$ , and let  $\mathcal{P}$  be the set of all measurable functions  $p : \mathbb{R}^d \to [0, \infty)$  such that  $\int p(\boldsymbol{x})(1 + \|\boldsymbol{x}\|^2) d\boldsymbol{x} < \infty$ . Now we want to maximize H(p) over all probability densities  $p \in \mathcal{P}$  satisfying

$$\int \boldsymbol{x} p(\boldsymbol{x}) \, d\boldsymbol{x} = \boldsymbol{\mu}_o \quad \text{and} \quad \int (\boldsymbol{x} - \boldsymbol{\mu}_o) (\boldsymbol{x} - \boldsymbol{\mu}_o)^\top p(\boldsymbol{x}) \, d\boldsymbol{x} = \boldsymbol{\Sigma}_o$$

for a given vector  $\boldsymbol{\mu}_o \in \mathbb{R}^d$  and a given symmetric and positive definite matrix  $\boldsymbol{\Sigma}_o \in \mathbb{R}^{d \times d}$ . Such a density p describes the distribution of a random vector  $\boldsymbol{X} \in \mathbb{R}^d$  with  $\mathbb{E}(\boldsymbol{X}) = \boldsymbol{\mu}_o$  and  $\operatorname{Var}(\boldsymbol{X}) = \boldsymbol{\Sigma}_o$ .

In other words, we want to minimize f(p) = -H(p) under the constraints that

$$egin{aligned} g_0(p) &:= \int p(oldsymbol{x}) \, doldsymbol{x} \ \stackrel{!}{=} \ 1, \ oldsymbol{g}_1(p) &:= \int oldsymbol{x} p(oldsymbol{x}) \, doldsymbol{x} \ \stackrel{!}{=} \ oldsymbol{\mu}_o, \ oldsymbol{G}_2(p) &:= \int (oldsymbol{x} - oldsymbol{\mu}_o) (oldsymbol{x} - oldsymbol{\mu}_o)^ op p(oldsymbol{x}) \, doldsymbol{x} \ \stackrel{l}{=} \ oldsymbol{\Sigma}_o. \end{aligned}$$

To this end we consider for a number  $\lambda_0 \in \mathbb{R}$ , a vector  $\lambda_1 \in \mathbb{R}^d$  and a symmetric matrix  $\Lambda_2 \in \mathbb{R}^{d \times d}$  the quantity

$$-H(p) - \lambda_0 g_0(p) - \boldsymbol{\lambda}_1^{\top} \boldsymbol{g}_1(p) - \operatorname{trace}(\boldsymbol{\Lambda}_2 \boldsymbol{G}_2(p)) \\ = \int \Big( \log p(\boldsymbol{x}) - \lambda_0 - \boldsymbol{\lambda}_1^{\top} \boldsymbol{x} - (\boldsymbol{x} - \boldsymbol{\mu}_o)^{\top} \boldsymbol{\Lambda}_2(\boldsymbol{x} - \boldsymbol{\mu}_o) \Big) p(\boldsymbol{x}) \, d\boldsymbol{x}.$$

The minimizer of that auxiliary functional is unique (almost everywhere) and given by

$$p(\boldsymbol{x}) = \exp\left(-1 + \lambda_0 + \boldsymbol{\lambda}_1^\top \boldsymbol{x} + (\boldsymbol{x} - \boldsymbol{\mu}_o)^\top \boldsymbol{\Lambda}_2(\boldsymbol{x} - \boldsymbol{\mu}_o)\right),$$

provided that  $-\Lambda_2$  is positive definite. Note that this is always a positive multiple of a Gaussian density. If we choose  $\Lambda_2 := -\Sigma_o^{-1}/2$ ,  $\lambda_1 := 0$  and the constant  $\lambda_0 := 1 - 2^{-1} d \log(2\pi) - 2^{-1} \log \det(\Sigma_o)$ , then

$$p(\boldsymbol{x}) = (2\pi)^{-d/2} \det(\boldsymbol{\Sigma}_o)^{-1/2} \exp\left(-(\boldsymbol{x}-\boldsymbol{\mu}_o)^\top \boldsymbol{\Sigma}_o^{-1} (\boldsymbol{x}-\boldsymbol{\mu}_o)/2\right),$$

which is the density function of the *d*-variate Gaussian distribution  $\mathcal{N}_d(\boldsymbol{\mu}_o, \boldsymbol{\Sigma}_o)$ .

This shows that the Gaussian distribution  $\mathcal{N}_d(\mu_o, \Sigma_o)$  is the most "diffuse" distribution (in the sense of entropy) among all probability distributions with given mean vector  $\mu_o$  and covariance matrix  $\Sigma_o$ .

**Exercise 5.7.** Determine a probability density p on  $(0, \infty)$  such that its entropy H(p) is maximized under the constraint

$$\int_0^\infty x p(x) \, dx = \mu > 0.$$

**Exercise 5.8** (Discrete Shannon entropy, I). For a vector  $\boldsymbol{p} \in [0,\infty)^d$  let

$$H(\boldsymbol{p}) := -\sum_{i=1}^d p_i \log p_i.$$

Maximize  $H(\mathbf{p})$  under the constraint that  $\sum_{i=1}^{d} p_i = 1$ .

**Exercise 5.9** (Discrete Shannon entropy, II). Let  $P \in [0, 1]^{r \times s}$  be a matrix with given row and column sums:

$$\sum_{j=1}^{\circ} P_{ij} = a_i \quad \text{for } i = 1, \dots, r$$

and

$$\sum_{i=1}^{r} P_{ij} = b_j \quad \text{for } j = 1, \dots, s$$

with given probability vectors  $\boldsymbol{a} \in (0,1)^r$  and  $\boldsymbol{b} \in (0,1)^s$ . Determine such a matrix  $\boldsymbol{P}$  with maximal entropy

$$H(\boldsymbol{P}) := -\sum_{i,j} P_{ij} \log P_{ij}.$$

Interpretation: The matrix P describes the joint distribution of two random variables X and Y with values in  $\{1, ..., r\}$  and  $\{1, ..., s\}$ , respectively. Their respective marginal distributions are given by a and b, and we are looking for a maximally diffuse joint distribution.

**Kullback–Leibler divergence.** Let P and Q be probability distributions on a measurable space  $(\mathcal{X}, \mathcal{B})$ . The Kullback–Leibler divergence of Q with respect to P is defined as

$$K(Q \mid P) := \int \log(q) \, dQ \in [0, \infty]$$

if Q has a density q with respect to P,<sup>1</sup> and  $K(Q | P) := \infty$  otherwise. Note that

$$\int \log(q) \, dQ = -\int_{\{q>0\}} \log(1/q)q \, dP \ge -\int_{\{q>0\}} (1/q-1)q \, dP = 1 - P\{q>0\} \ge 0$$

 ${}^{1}Q(B) = \int_{B} q \, dP$  for arbitrary  $B \in \mathcal{B}$ 

with equality if and only if q = 1 *P*-almost everywhere. Thus,  $K(Q | P) \ge 0$  with equality if and only if  $Q \equiv P$ .

By means of Lagrange's method one can show that

$$K(Q \mid P) = \sup_{g \in \mathcal{G}(P)} \int g \, dQ$$

where  $\mathcal{G}(P)$  is the set of all measurable functions  $g : \mathcal{X} \to [-\infty, \infty)$  such that  $\int g \, dQ$  exists in  $[-\infty, \infty]$  and

$$\int e^g \, dP = 1.$$

Suppose first that Q(B) > 0 = P(B) for some  $B \in \mathcal{B}$ . Then for r > 0,  $g_r := 1_B \cdot r$  defines a function in  $\mathcal{G}(P)$  such that  $\int g_r dQ = rQ(B) \to \infty$  as  $r \to \infty$ . Note also that the existence of such a set B implies that Q does not admit a density with respect to P, whence  $K(Q | P) = \infty$ .

If no such set *B* exists, it follows from the theorem of Radon–Nikodym that *Q* does admit a density *q* with respect to *P*. Then we may write  $\int g \, dQ = \int gq \, dP$ , and we consider the larger family  $\overline{\mathcal{G}}(P)$  of all measurable functions  $g : \mathcal{X} \to [-\infty, \infty)$  such that  $\int g \, dQ$  exists in  $[-\infty, \infty]$  and  $\int e^g \, dP < \infty$ . Instead of maximising  $\int g \, dQ$  over all  $g \in \mathcal{G}(P)$ , we try to maximise

$$\int g \, dQ - \int e^g \, dP = \int (gq - e^g) \, dF$$

over all functions  $g \in \overline{\mathcal{G}}(P)$ , which is Lagrange's method with  $\lambda = 1$ . Indeed, for any  $x \in \mathcal{X}$ ,

$$\underset{h \in [-\infty,\infty)}{\operatorname{arg\,max}} (hq(x) - e^h) = \log q(x)$$

(with  $\log 0 := -\infty$ ), so

$$\int (gq - e^g) dP \leq \int (\log(q) - 1)q dP = \int \log(q) dQ - 1$$

with equality if and only if  $g = g_* := \log(q)$  *P*-almost everywhere. Since  $g_* \in \mathcal{G}(P)$ , this shows that  $K(Q | P) = \int g_* dQ = \max_{g \in \mathcal{G}(P)} \int g dQ$ .

**Exercise 5.10.** Given a strictly positive probability density  $p_o$  on (0, 1) and numbers  $\mu_o, \sigma_o \in (0, 1)$ , determine a probability density p with minimal Kullback–Leibler divergence  $K(p | p_o) := \int_0^1 \log(p(x)/p_o(x))p(x) dx$  under the constraints that

$$\mu(p) \ := \ \int_0^1 x p(x) \, dx \ = \ \mu_o \quad \text{and} \quad \sigma^2(p) \ := \ \int_0^1 (x - \mu_o)^2 p(x) \, dx \ = \ \sigma_o^2.$$

With a suitable Lagrange ansatz you'll see immediately how the density p should look like.

(For which values of  $\mu_o$  and  $\sigma_o$  does a solution exist?)

**Optimal kernel functions.** In connection with kernel density estimation as well as in connection with Wilcoxon's rank tests one encounters the following optimization problem: Find a measurable

function  $K:\mathbb{R}\rightarrow [0,\infty)$  such that

$$g_0(K) := \int K(x) \, dx \stackrel{!}{=} 1,$$
  
$$g_1(K) := \int K(x) x \, dx \stackrel{!}{=} 0,$$
  
$$g_2(K) := \int K(x) x^2 \, dx \stackrel{!}{=} 1$$

with minimal value of

$$f(K) := \int K(x)^2 \, dx.$$

To solve this problem, let  $\mathcal{K}$  be the set of all measurable functions  $K : \mathbb{R} \to [0, \infty)$  such that

$$\int K(x)(1+x^2)\,dx < \infty.$$

We just ignore the constraint  $g_1(K) = 0$ , hoping that the minimizer  $K \in \mathcal{K}$  of f(K) under the constraints  $g_0(K) = g_2(K) = 1$  will automatically satisfy  $g_1(K) = 0$ .

The Lagrange ansatz leads to the auxiliary function

$$f(K) - \lambda_0 g_0(K) - \lambda_2 g_2(K) = \int (K(x)^2 - K(x)(\lambda_0 + \lambda_2 x^2)) dx$$

with a certain vector  $\boldsymbol{\lambda} = (\lambda_0, \lambda_2)^\top \in \mathbb{R}^2$ . Now

$$K(x)^{2} - K(x)(\lambda_{0} + \lambda_{2}x^{2}) = (K(x) - (\lambda_{0} + \lambda_{2}x^{2})/2)^{2} - 4^{-1}(\lambda_{0} + \lambda_{2}x^{2})^{2},$$

so an optimal function K would be given by

$$K_{\lambda}(x) = (\lambda_0/2 + \lambda_2 x^2/2)^+.$$

This defines a nontrivial function in  $\mathcal{K}$  if and only if  $\lambda_0 > 0 > \lambda_2$ . With constants a, b > 0 we may also consider

$$K_{a,b}(x) := a(1 - (x/b)^2)^+.$$

Note first that with the substition z = x/b, dx = b dz,

$$g_0(K_{a,b}) = a \int_{-b}^{b} (1 - (x/b)^2) dx$$
  
=  $ab \int_{-1}^{1} (1 - z^2) dz = ab \cdot 2(z - z^3/3) \Big|_{z=0}^{1} = \frac{4ab}{3},$ 

so

$$a \stackrel{!}{=} \frac{3}{4b}.$$

Furthermore, since  $K_{a,b}$  is an even function,  $g_1(K_{a,b}) = 0$ . Finally,

$$g_2(K_{a,b}) = a \int_{-b}^{b} (1 - (x/b)^2) x^2 dx$$
  
=  $\frac{3b^2}{4} \int_{-1}^{1} (1 - z^2) z^2 dz = ab^3 \cdot 2(z^3/3 - z^5/5) \Big|_{z=0}^{1} = \frac{b^2}{5},$ 

so

$$b \stackrel{!}{=} \sqrt{5}$$

All in all, the original optimization problem has the (almost everywhere) unique solution

$$K(x) = \frac{3}{4\sqrt{5}} \left(1 - \frac{x^2}{5}\right)^+,$$

a so-called Epanechnikov kernel.

The shape of a hanging chain. Let us consider a thin chain of length 2L which is hanging freely from two points at identical height and with distance 2M, where 0 < M < L. The question is which shape of the chain will materialize. We talk about chains rather than strings or ropes to emphasize that within the chain there are no elastic forces at work. Hence the only driving forces result from gravitation.

Formulating the problem mathematically, we are looking for a curve

$$[-L,L] \ni t \mapsto [x(t),y(t)]^{\top} \in \mathbb{R}^2$$

with the following properties:

• The functions x and y are continuously differentiable<sup>2</sup> on [-L, L] with derivatives x' and y' satisfying

$$x'(t)^2 + y'(t)^2 = 1$$

for arbitrary  $t \in [-L, L]$ . That means, our curve is parametrized in arc length.

• The curve starts in  $(-M, 0)^{\top}$  and ends in  $(M, 0)^{\top}$ , that means,

$$x(\pm L) = \pm M$$
 and  $y(\pm L) = 0$ .

Under these constraints we would like to minimize the potential energy

$$E := \int_{-L}^{L} y(t) dt.$$

This definition comes from the consideration that for  $t \in [-L, L)$  and a small  $\delta > 0$ , the section  $\{(x(s), y(s)) : t \le s < t + \delta\}$  of our chain has mass proportional to  $\delta$  and is hanging at height y(t).

To solve this minimization problem we introduce

$$\boldsymbol{v}(t) := \left[ x'(t), y'(t) \right]^{\top},$$

defining a measurable mapping v from [-L, L] to the unit sphere of  $\mathbb{R}^2$ . Now we rewrite E and the constraints in terms of v: On the one hand, it follows from y(-L) = 0 and y(L) - y(-L) = 0

<sup>&</sup>lt;sup>2</sup>it would even suffice to assume that x and y are absolutely continuous.

 $\int_{-L}^{L} y'(t) dt = 0$  that

$$E(\boldsymbol{v}) = \int_{-L}^{L} \int_{-L}^{t} y'(s) \, ds \, dt$$
$$= \int_{-L}^{L} \int_{s}^{L} y'(s) \, dt \, ds$$
$$= \int_{-L}^{L} y'(s)(L-s) \, ds$$
$$= \int_{-L}^{L} y'(s)(-s) \, ds.$$

On the other hand, the constraint x(L) - x(-L) = 2M is equivalent to

$$G(\boldsymbol{v}) := \int_{-L}^{L} x'(s) \, ds = 2M.$$

Consequently we are now trying the following ansatz: For some  $\lambda \in \mathbb{R}$  we want to minimize the function

$$E(\boldsymbol{v}) - \lambda G(\boldsymbol{v}) = \int_{-L}^{L} \left( -y'(s)s - \lambda x'(s) \right) ds = -\int_{-L}^{L} \left\langle \boldsymbol{v}(s), (\lambda, s)^{\top} \right\rangle ds$$

with respect to v. The solution may be determined point-wise: By means of the Cauchy–Schwarz inequality, the optimal v is given by

$$\boldsymbol{v}_{\lambda}(s) := \left(\frac{\lambda}{\sqrt{\lambda^2 + s^2}}, \frac{s}{\sqrt{\lambda^2 + s^2}}\right)^{\top},$$

obviously a continuous function of  $s \in [-L, L]$ . Only  $\lambda > 0$  yields a reasonable solution because for  $\lambda \leq 0$  we would have  $G(\boldsymbol{v}_{\lambda}) \leq 0$ . For the curve (x, y) itself this yields the following solution: Note that

$$\int \frac{s}{\sqrt{\lambda^2 + s^2}} \, ds = C + \sqrt{\lambda^2 + s^2}.$$

Combining this with the boundary conditions that  $y_{\lambda}(\pm L) = 0$ , we obtain

$$y_{\lambda}(s) = \sqrt{\lambda^2 + s^2} - \sqrt{\lambda^2 + L^2}.$$

Recall the hyperbolic functions  $\cosh, \sinh : \mathbb{R} \to \mathbb{R}$  with  $\cosh(r) = (e^r + e^{-r})/2$  and  $\sinh(r) = (e^r - e^{-r})/2$ . They satisfy the equation

$$\cosh = \sqrt{1 + \sinh^2},$$

while their derivatives are given by  $\cosh' = \sinh \text{ and } \sinh' = \cosh > 0$ . Moreover,  $\sinh \text{ is}$  bijective with inverse function  $\operatorname{arsinh} : \mathbb{R} \to \mathbb{R}$  such that

$$\operatorname{arsinh}'(r) = \frac{1}{\sqrt{1+r^2}}$$

for arbitrary  $r \in \mathbb{R}$ . From this and the boundary conditions  $x_{\lambda}(\pm L) = \pm M$  one can easily derive that

$$x_{\lambda}(s) = \lambda \operatorname{arsinh}(s/\lambda),$$

118



Figure 5.1: Some catenaries.

and we may rewrite this as

$$s = \lambda \sinh(x_{\lambda}(s)/\lambda).$$

Now the boundary condition  $x_{\lambda}(\pm L) = \pm M$  is equivalent to the requirement that  $L/\lambda = \sinh(M/\lambda)$ , that is,

(5.3) 
$$L/M = \frac{\sinh(M/\lambda)}{M/\lambda} = 1 + \sum_{k=1}^{\infty} \frac{(M/\lambda)^{2k}}{(2k+1)!}$$

Since the right hand side is continuous and strictly decreasing in  $\lambda > 0$  with limits  $\infty$  as  $\lambda \to 0$ and 1 as  $\lambda \to \infty$ , there exists a unique solution  $\lambda > 0$  of the equation (5.3). For this particular  $\lambda > 0$ , the functions  $x(\cdot) := x_{\lambda}(\cdot)$  and  $y(\cdot) := y_{\lambda}(\cdot)$  are the desired solutions with

$$y = \lambda \sqrt{1 + \sinh(x/\lambda)^2} - \lambda \sqrt{1 + (L/\lambda)^2}$$
  
=  $\lambda \sqrt{1 + \sinh(x/\lambda)^2} - \lambda \sqrt{1 + \sinh(M/\lambda)^2}$   
=  $\lambda \cosh(x/\lambda) - \lambda \cosh(M/\lambda).$ 

Hence the chain configuration is given by the set

$$\Big\{ \big(x, \,\lambda \cosh(x/\lambda) - \lambda \cosh(M/\lambda)\big)^\top : x \in [-M, M] \Big\}.$$

Such a curve is also called a *catenary*. One of the first people to propose catenaries in architecture was the English natural philosopher, architect and polymath Robert Hooke (1635–1703). The famous Catalan architect Antoni Gaudi (1852–1926) used chain models to design a (never completed) church at Santa Coloma de Cervello. Figure 5.1 shows three different catenaries.

**Exercise 5.11** (A smooth ride with the subway). (a) Determine for a given parameter c > 0 a twice continuously differentiable function  $h : [0, 1] \to \mathbb{R}$  such that

$$h(0) = 0, \quad h(1) = 1, \quad h'(0) = h'(1) = 0, \quad |h''| \le c$$

with minimal value of

$$U(h) := \int_0^1 h''(t)^2 dt$$

For which parameters c does a solution exist at all?

(b) One may view h(t) as the position of a subway along a given route at time t. At time 0 the subway starts at station A, and at time 1 it should arrive at station B. The cost functional U(h) quantifies how problematic the ride would be for the material. For passengers, however, sudden changes in h'' or values  $h''(0), h''(1) \neq 0$  would be considered as rather unpleasant. How would the solution to (a) change if we impose that h'' should be Lipschitz-continuous with constant  $c_3 > 0$  and satisfy h''(0) = h''(1) = 0?

**The Neyman–Pearson lemma.** Suppose we observe a random variable X with unknown distribution P on a measurable space  $(\mathcal{X}, \mathcal{B})$ . We would like to test the null hypothesis " $P = P_0$ " versus the alternative hypothesis " $P = P_1$ " at a given test level  $\alpha \in (0, 1)$ , where  $P_0, P_1$  are given probability distributions on  $(\mathcal{X}, \mathcal{B})$ .

That means, we are looking for a statistical test  $\phi$ , i.e. a measurable function  $\phi : \mathcal{X} \to [0, 1]$ , where we reject " $P = P_0$ " with probability  $\phi(X)$ . The probability of an error of the first kind equals

$$g(\phi) := \int_{\mathcal{X}} \phi \, dP_0$$

and should not exceed  $\alpha$ , while the probability of an error of the second kind equals

$$f(\phi) := 1 - \int_{\mathcal{X}} \phi \, dP_1$$

and should be as small as possible.

To solve this optimization problem, let  $P_0$  and  $P_1$  have densities  $p_0$  and  $p_1$ , respectively, with respect to some measure  $\mu$  on  $(\mathcal{X}, \mathcal{B})$ . This is not a restriction: One could choose  $\mu = P_0 + P_1$ and deduce from the Radon–Nikodym theorem the existence of such densities  $p_0$  and  $p_1$ . Now we make the Lagrange ansatz and consider for  $\lambda \in \mathbb{R}$  the functional

$$\phi \mapsto f(\phi) + \lambda g(\phi) = 1 + \int_{\mathcal{X}} \phi(\lambda p_0 - p_1) d\mu.$$

(Using  $+\lambda g(\phi)$  instead of  $-\lambda g(\phi)$  will be convenient notationally.) This functional may be minimized by minimizing the integrand pointwise: A test  $\phi$  minimizes  $f + \lambda g$  if and only if

$$\phi = \begin{cases} 1 & \mu\text{-almost everywhere on } \{p_1 > \lambda p_0\}, \\ 0 & \mu\text{-almost everywhere on } \{p_1 < \lambda p_0\}. \end{cases}$$

121

Specifically, let

$$\phi_{\lambda,\gamma}(x) := \begin{cases} 1 & \text{if } p_1(x) > \lambda p_0(x) \\ \gamma & \text{if } p_1(x) = \lambda p_0(x) \\ 0 & \text{if } p_1(x) < \lambda p_0(x) \end{cases}$$

with  $\gamma \in [0, 1]$ . Then

$$\int_{\mathcal{X}} \phi_{\lambda,\gamma} \, dP_0 = P_0\{p_1 > \lambda p_0\} + \gamma P_0\{p_1 = \lambda p_0\}$$

Note that  $H(\lambda) := P_0\{p_1 > \lambda p_0\}$  is monotone decreasing and right-continuous in  $\lambda \in \mathbb{R}$  with

$$H(\lambda -) - H(\lambda) = P_0\{p_1 = \lambda p_0\},$$

so

$$\int_{\mathcal{X}} \phi_{\lambda,\gamma} \, dP_0 = H(\lambda) + \gamma \big( H(\lambda) - H(\lambda -) \big).$$

Moreover, H(0-) = 1 and  $\lim_{\lambda \to \infty} H(\lambda) = 0$ . Hence an optimal test of " $P = P_0$ " versus " $P = P_1$ " at level  $\alpha$  is given by  $\phi_{\lambda,\gamma}$ , where

$$\begin{split} \lambda &:= \min\{\lambda \in \mathbb{R} : H(\lambda) \leq \alpha\} \geq 0, \\ \gamma &:= \begin{cases} 0 & \text{if } H(\lambda) = \alpha, \\ \frac{\alpha - H(\lambda)}{H(\lambda -) - H(\lambda)} & \text{if } H(\lambda) > \alpha. \end{cases} \end{split}$$

This optimal test was developed by Neyman<sup>3</sup> and Pearson<sup>4</sup> (1933).

# 5.1.3 Justification of Lagrange's method

At first glance, Lagrange's method looks like a cute trick which may work or fail, but we don't know when and why. By means of our results from convex analysis we can show that for convex functions  $f : \mathbb{R}^d \to \mathbb{R}$  and linear functions  $g : \mathbb{R}^d \to \mathbb{R}^q$  the method has to work.

**Theorem 5.12**. Let  $\mathcal{X}$  be an open convex subset of  $\mathbb{R}^d$ , let  $f : \mathcal{X} \to \mathbb{R}$  be a convex and  $g : \mathbb{R}^d \to \mathbb{R}^q$  be a linear function. Furthermore let C be a convex subset of  $\mathbb{R}^q$ . Suppose there exists a point

$$oldsymbol{x}_* \in rgmin_{oldsymbol{x}\in\mathcal{X}\,:\,oldsymbol{g}(oldsymbol{x})\inoldsymbol{C}} f(oldsymbol{x}).$$

Then there exists a vector  $\boldsymbol{\lambda} \in \mathbb{R}^q$  such that

$$oldsymbol{x}_{*} \in rgmin_{oldsymbol{x}\in\mathcal{X}}ig(f(oldsymbol{x})-oldsymbol{\lambda}^{ op}oldsymbol{g}(oldsymbol{x})ig)$$

and

$$oldsymbol{\lambda}^ opoldsymbol{g}(oldsymbol{x}) \, \geq \, oldsymbol{\lambda}^ opoldsymbol{g}(oldsymbol{x}_*) \quad ext{whenever} \, oldsymbol{x} \in \mathcal{X}, oldsymbol{g}(oldsymbol{x}) \in oldsymbol{C}.$$

<sup>&</sup>lt;sup>3</sup>Jerzy Neyman (1894–1981): Polish mathematician and statistician; pioneered hypothesis testing and confidence intervals.

<sup>&</sup>lt;sup>4</sup>Egon S. Pearson (1895–1980): British statistician; son of the mathematician and biostatistician Karl Pearson (1857–1936).

Finally we formulate an important special case of Theorems 5.1 and 5.12 concerning optimization problems with linear inequality and equality constraints:

**Theorem 5.13** (Karush–Kuhn–Tucker conditions). Let  $\mathcal{X} \subset \mathbb{R}^d$  be a convex set,  $f : \mathcal{X} \to \mathbb{R}$  a convex function and  $g = (g_j)_{j=1}^q : \mathbb{R}^d \to \mathbb{R}^q$  a linear function. For certain numbers  $c_1, c_2, \ldots, c_q \in \mathbb{R}$  and  $q_o \in \{0, 1, \ldots, q\}$  let

$$\mathcal{K} := \left\{ \boldsymbol{x} \in \mathcal{X} : g_j(\boldsymbol{x}) \le c_j \text{ if } j \le q_o, g_j(\boldsymbol{x}) = c_j \text{ if } j > q_o \right\}.$$

(a) Suppose there exist  $\lambda \in \mathbb{R}^q$  and

$$oldsymbol{x}_* \ \in \ rgmin_{oldsymbol{x} \in \mathcal{X}} \min ig(f(oldsymbol{x}) - oldsymbol{\lambda}^ op oldsymbol{g}(oldsymbol{x})ig)$$

such that  $x_* \in \mathcal{K}$  and

(5.4) 
$$\lambda_j \begin{cases} \leq 0 & \text{if } j \leq q_o, \\ = 0 & \text{if } j \leq q_o \text{ and } g_j(\boldsymbol{x}_*) < c_j. \end{cases}$$

Then

$$oldsymbol{x}_* \in \operatorname*{arg\,min}_{oldsymbol{x}\in\mathcal{K}} f(oldsymbol{x}).$$

(**b**) Suppose that  $\mathcal{X}$  is open, and let

$$oldsymbol{x}_* \in rgmin_{oldsymbol{x}\in\mathcal{K}} f(oldsymbol{x}).$$

Then there exists a vector  $\boldsymbol{\lambda} \in \mathbb{R}^q$  such that

$$oldsymbol{x}_{*} \ \in \ rgmin_{oldsymbol{x} \in \mathcal{X}} \min ig(f(oldsymbol{x}) - oldsymbol{\lambda}^{ op}oldsymbol{g}(oldsymbol{x})ig)$$

and (5.4) is satisfied.

**Proof of Theorem 5.12.** Let  $\mathcal{K} = g^{-1}(C) \cap \mathcal{X}$ . By linearity of g, convexity of C and  $\mathcal{X}$  implies that  $\mathcal{K}$  is a convex set, too. Thus

$$\mathcal{D} := \mathcal{K} imes (-\infty, f(\boldsymbol{x}_*)) \quad ext{and} \quad ext{epi}(f) = \left\{ (\boldsymbol{x}, r) \in \mathcal{X} imes \mathbb{R} : f(\boldsymbol{x}) \leq r 
ight\}$$

are convex and disjoint subsets of  $\mathbb{R}^d \times \mathbb{R}$ . According to Theorem 2.22, there exists a nonzero pair  $(\mathbf{b}, t) \in \mathbb{R}^d \times \mathbb{R}$  such that

$$\boldsymbol{b}^{\top}\boldsymbol{y} + ts \leq \boldsymbol{b}^{\top}\boldsymbol{x} + tr \text{ for all } (\boldsymbol{y}, s) \in \mathcal{D}, (\boldsymbol{x}, r) \in \operatorname{epi}(f).$$

Setting  $y = x = x_*$  and  $s = f(x_*) - 1$ ,  $r = f(x_*)$  shows that  $t \ge 0$ . Moreover, t = 0 would imply that  $b \ne 0$  and  $b^{\top} x \ge b^{\top} x_*$  for all  $x \in \mathcal{X}$ , a contradiction to  $\mathcal{X}$  being an open set and  $x_* \in \mathcal{X}$ . Hence we may assume without loss of generality that t = 1. Then the separation inequality is equivalent to

(5.5) 
$$\boldsymbol{b}^{\top}\boldsymbol{y} + f(\boldsymbol{x}_*) \leq \boldsymbol{b}^{\top}\boldsymbol{x} + f(\boldsymbol{x}) \text{ for all } \boldsymbol{y} \in \mathcal{K}, \boldsymbol{x} \in \mathcal{X}.$$

Setting  $x = x_*$  in (5.5) leads to

(5.6) 
$$\boldsymbol{b}^{\top} \boldsymbol{y} \leq \boldsymbol{b}^{\top} \boldsymbol{x}_{*}$$
 for all  $\boldsymbol{y} \in \mathcal{K}$ 

But for sufficiently small  $\delta > 0$ , all points  $\boldsymbol{y} = \boldsymbol{x}_* \pm \boldsymbol{w}$  with  $\boldsymbol{w} \in \mathbb{R}^d$ ,  $\boldsymbol{g}(\boldsymbol{w}) = \boldsymbol{0}$  and  $\|\boldsymbol{w}\| < \delta$  belong to the set  $\mathcal{K}$ . Hence it follows from (5.6) that  $\boldsymbol{b}^\top \boldsymbol{w} = 0$  for all  $\boldsymbol{w} \in \mathbb{R}^d$  such that  $\boldsymbol{g}(\boldsymbol{w}) = \boldsymbol{0}$ . If we write  $\boldsymbol{g}(\boldsymbol{w}) = \boldsymbol{B}^\top \boldsymbol{w}$  with a matrix

$$oldsymbol{B} = [oldsymbol{b}_1, oldsymbol{b}_2, \dots, oldsymbol{b}_q] \in \mathbb{R}^{d imes q},$$

then **b** is perpendicular to all vectors in  $\{\boldsymbol{b}_1, \ldots, \boldsymbol{b}_q\}^{\perp}$ . In other words, **b** has to be a linear combination of the vectors  $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_q$ , see Exercise 5.14. Hence  $\boldsymbol{b} = -\sum_{j=1}^q \lambda_j \boldsymbol{b}_j = -\boldsymbol{B}\boldsymbol{\lambda}$  for some  $\boldsymbol{\lambda} \in \mathbb{R}^q$ , so  $\boldsymbol{b}^\top \boldsymbol{x} = -\boldsymbol{\lambda}^\top \boldsymbol{B}^\top \boldsymbol{x} = -\boldsymbol{\lambda}^\top \boldsymbol{g}(\boldsymbol{x})$ . Consequently we may deduce from (5.5) that

$$f(\boldsymbol{x}) - \boldsymbol{\lambda}^{ op} \boldsymbol{g}(\boldsymbol{x}) \geq f(\boldsymbol{x}_*) - \boldsymbol{\lambda}^{ op} \boldsymbol{g}(\boldsymbol{x}_*) \quad ext{for all } \boldsymbol{x} \in \mathcal{X}$$

while (5.6) is equivalent to

$$oldsymbol{\lambda}^ op g(oldsymbol{y}) \, \geq \, oldsymbol{\lambda}^ op g(oldsymbol{x}_*) \quad ext{for all } oldsymbol{y} \in \mathcal{K}.$$

**Exercise 5.14** (Bi-orthogonal sets). Show that for any nonvoid set  $M \subset \mathbb{R}^d$ ,

$$(\boldsymbol{M}^{\perp})^{\perp} = \operatorname{span}(\boldsymbol{M}),$$

where  $M^{\perp} := \{ \boldsymbol{y} \in \mathbb{R}^d : \boldsymbol{x}^{\top} \boldsymbol{y} = 0 \text{ for all } \boldsymbol{x} \in \boldsymbol{M} \}.$ 

**Proof of Theorem 5.13.** Part (a) follows immediately from Theorem 5.1 and Remarks 5.2 and 5.3.

It remains to prove part (b). According to Theorem 5.12 there exists a vector  $\lambda \in \mathbb{R}^q$  such that  $x_*$  minimizes  $f - \lambda^\top g$  on  $\mathcal{X}$  and

$$oldsymbol{\lambda}^ op oldsymbol{g}(oldsymbol{y}) \ \geq oldsymbol{\lambda}^ op oldsymbol{g}(oldsymbol{x}_*) \quad ext{for all }oldsymbol{y} \in \mathcal{K}.$$

In other words,  $\lambda^{\top} g(v) \ge 0$  for any vector  $v \in \mathbb{R}^q$  such that  $x_* + tv \in \mathcal{K}$  for some t > 0. Since  $x_*$  is an interior point of  $\mathcal{X}$ , one may reformulate this statement as follows: For  $v \in \mathbb{R}^d$ ,

(5.7) 
$$\boldsymbol{\lambda}^{\top} \boldsymbol{g}(\boldsymbol{v}) \geq 0$$
 whenever  $g_i(\boldsymbol{v}) \begin{cases} \leq 0 & \text{if } i \leq q_o \text{ and } g_i(\boldsymbol{x}_*) = c_i, \\ = 0 & \text{if } i > q_o. \end{cases}$ 

Unfortunately, the latter condition on  $\lambda$  alone does not guarantee the KKT conditions (5.4). To achieve the latter, we have to "sparsify"  $\lambda$  as follows: We write  $g_j(v) = b_j^\top v$  with certain vectors  $b_1, \ldots, b_q \in \mathbb{R}^d$  and define the index set

$$J = J(\boldsymbol{\lambda}) := \{ j \in \{1, \dots, q\} : \lambda_j \neq 0 \text{ or } g_j(\boldsymbol{x}_*) = c_j \}.$$

Suppose that for some index j,

$$\lambda_j \neq 0$$
 and  $\boldsymbol{b}_j \in \operatorname{span}(\boldsymbol{b}_i : i \in J \setminus \{j\}).$ 

Then there exist real numbers  $\mu_i, i \in J \setminus \{j\}$  such that  $b_j = \sum_{i \in J \setminus \{j\}} \mu_i b_i$ . But then

$$oldsymbol{\lambda}^ op g \,\equiv\, { ilde\lambda}^ op g$$

with

$$ilde{\lambda}_i := egin{cases} \lambda_i + \mu_i \lambda_j & ext{if } i \in J \setminus \{j\}, \ 0 & ext{else.} \end{cases}$$

Hence we may modify the vector  $\lambda$  step by step without changing  $\lambda^{\top} g$  until finally its index set J has the following property:

(5.8) 
$$\mathbf{b}_j \notin \operatorname{span}(\mathbf{b}_i : i \in J \setminus \{j\})$$
 whenever  $\lambda_j \neq 0$ .

This vector  $\lambda$  does satisfy the KKT conditions (5.4). Suppose the contrary. Then there exists an index  $j \leq q_o$  such that either  $\lambda_j > 0$  or  $\lambda_j < 0$  and  $g_j(\boldsymbol{x}_*) < c_j$ . By property (5.8) there exists a vector  $\boldsymbol{w} \in \mathbb{R}^d$  such that

$$\boldsymbol{b}_j^\top \boldsymbol{w} = 1$$
 and  $\boldsymbol{b}_i^\top \boldsymbol{w} = 0$  for all  $i \in J \setminus \{j\}$ .

But then the vector  $\boldsymbol{v} := -|\lambda_j|\boldsymbol{w}$  would satisfy the inequalities for  $g_i(\boldsymbol{v})$  in (5.7), because  $g_i(\boldsymbol{v}) = -|\lambda_j|\mathbf{1}_{[i=j]}$  for  $i \in J$ . But  $\boldsymbol{\lambda}^\top \boldsymbol{g}(\boldsymbol{v}) = -\lambda_j^2 < 0$ . This contradiction shows that  $\boldsymbol{\lambda}$  has to satisfy (5.4).

**Exercise 5.15** (Empirical likelihood). Let  $x_1, \ldots, x_n$  be pairwise different vectors in  $\mathbb{R}^d$  which are not lying on a common hyperplane. Show that for any interior point  $\mu$  of  $conv(x_1, \ldots, x_n)$  there exists at least one vector  $p \in (0, 1)^n$  such that

$$\sum_{i=1}^n p_i = 1$$
 and  $\sum_{i=1}^n p_i x_i = \mu$ .

(a) Show that among all such vectors p there is precisely one minimizer of

$$f(\boldsymbol{p}) := -\sum_{i=1}^n \log p_i.$$

Show that this minimizer has the form

$$p_i = (a + \boldsymbol{b}^\top \boldsymbol{x}_i)^{-1}$$

for suitable parameters  $a \in \mathbb{R}$  and  $b \in \mathbb{R}^d$ .

(b) Show that a vector  $p \in (0,1)^n$  solves the optimization problem in part (a) if and only if for some vector  $b \in \mathbb{R}^d$ ,

$$p_i = \frac{1}{n + \boldsymbol{b}^{\top}(\boldsymbol{x}_i - \boldsymbol{\mu})}$$
 with  $n + \boldsymbol{b}^{\top}(\boldsymbol{x}_i - \boldsymbol{\mu}) > 0$  for  $1 \le i \le n$ 

and

$$\sum_{i=1}^n rac{oldsymbol{x}_i - oldsymbol{\mu}}{n + oldsymbol{b}^ op (oldsymbol{x}_i - oldsymbol{\mu})} \ = \ oldsymbol{0}.$$

(c) Show that the vector **b** in part (b) is the unique minimizer of

$$\widetilde{f}(oldsymbol{b}) \ := \ -\sum_{i=1}^n \log ig(n+oldsymbol{b}^ op(oldsymbol{x}_i-oldsymbol{\mu})ig)$$

on  $\mathbb{R}^d$ , where  $\log(r) := -\infty$  for  $r \leq 0$ .

(d) Describe and implement a procedure to solve the optimization problem in part (a) by means of parts (b-c).

# 5.1.4 Lagrange duality

In connection with Theorem 5.12, there is a so-called dual approach to optimization under linear constraints. We elaborate on this in the following setting: Let  $f : \mathbb{R}^d \to \mathbb{R}$  be convex and coercive, and let  $g : \mathbb{R}^d \to \mathbb{R}^q$  be linear and surjective with q < d. The goal is to minimize f over  $N(g) := \{x \in \mathbb{R}^d : g(x) = 0\}.$ 

By coercivity of f, there exists a point

$$\boldsymbol{x}_* \in \operatorname*{arg\,min}_{\boldsymbol{x}\in N(\boldsymbol{g})} f(\boldsymbol{x}),$$

and by Theorem 5.12, there exists a vector  $\boldsymbol{\lambda}_* \in \mathbb{R}^q$  such that

$$oldsymbol{x}_* \in rgmin_{oldsymbol{x} \in \mathbb{R}^d}ig(f(oldsymbol{x}) - oldsymbol{\lambda}_*^ op oldsymbol{g}(oldsymbol{x})ig)$$

We define an auxiliary function  $L: \mathbb{R}^q \to [-\infty, \infty)$  via

$$L(oldsymbol{\lambda}) \ := \ \inf_{oldsymbol{x} \in \mathbb{R}^d} ig(f(oldsymbol{x}) - oldsymbol{\lambda}^ op oldsymbol{g}(oldsymbol{x})ig).$$

This function L is concave and upper semicontinuous, that is, -L is convex and lower semicontinuous. Indeed, if g(x) = Ax,  $x \in \mathbb{R}^d$ , for some matrix  $A \in \mathbb{R}^{q \times d}$ , then

$$L(\boldsymbol{\lambda}) = -f^*(\boldsymbol{A}^\top \boldsymbol{\lambda})$$

with the Fenchel-Laplace transform  $f^*$  of f; see Chapter 3. An important fact is that

$$f(\boldsymbol{x}_*) = L(\boldsymbol{\lambda}_*) = \max_{\boldsymbol{\lambda} \in \mathbb{R}^q} L(\boldsymbol{\lambda}).$$

Indeed, since  $\boldsymbol{x}_* \in N(\boldsymbol{g})$ ,

$$L(\boldsymbol{\lambda}_*) = f(\boldsymbol{x}_*) - \boldsymbol{\lambda}_*^\top \boldsymbol{g}(\boldsymbol{x}_*) = f(\boldsymbol{x}_*),$$

and for any  $\lambda \in \mathbb{R}^q$ ,

$$L(\boldsymbol{\lambda}) \leq f(\boldsymbol{x}_*) - \boldsymbol{\lambda}^{ op} \boldsymbol{g}(\boldsymbol{x}_*) = f(\boldsymbol{x}_*).$$

These findings suggest the following strategy to minimize f over N(g):

Step 1. Find a vector

$$\widehat{\boldsymbol{\lambda}} \in rg \max_{\boldsymbol{\lambda} \in \mathbb{R}^q} L(\boldsymbol{\lambda}).$$

Step 2. Determine a point

$$egin{aligned} \widehat{oldsymbol{x}} &\in rgmin_{oldsymbol{x} \in \mathbb{R}^d} ig( f(oldsymbol{x}) - \widehat{oldsymbol{\lambda}}^{ op} oldsymbol{g}(oldsymbol{x}) ig) \end{aligned}$$

and check whether  $\hat{x} \in N(g)$ . If yes, this point  $\hat{x}$  minimizes f over N(g).

**Remark 5.16**. If *L* is strictly concave on  $\{L > -\infty\}$ , then Step 1 yields automatically  $\hat{\lambda} = \lambda_*$ , and Step 2 is possible in the sense that some minimizer of  $f(x) - \hat{\lambda}^\top g(x)$  over all  $x \in \mathbb{R}^q$  does belong to N(g). More generally, let  $g^*$  be the conjugate mapping of g, that is, the mapping  $\lambda \mapsto A^\top \lambda$  if g is given by  $x \mapsto Ax$ . Since

$$f(\boldsymbol{x}) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(\boldsymbol{x}) = f(\boldsymbol{x}) - \boldsymbol{g}^{*}(\boldsymbol{\lambda})^{\top} \boldsymbol{x},$$

one can write

$$L(\boldsymbol{\lambda}) = \tilde{L}(\boldsymbol{g}^*(\boldsymbol{\lambda}))$$

with a concave function  $\tilde{L}: N(\boldsymbol{g})^{\perp} \to [-\infty, \infty)$ . Here we utilise the fact that  $N(\boldsymbol{g})^{\perp} = \boldsymbol{g}^*(\mathbb{R}^q)$ . If the latter function is strictly concave on  $\{\tilde{L} > -\infty\}$ , the set  $\hat{\mathcal{X}}$  coincides with the set  $\mathcal{X}_* = \arg \min_{\boldsymbol{x} \in \mathbb{R}^q} (f(\boldsymbol{x}) - \boldsymbol{\lambda}_*^{\top} \boldsymbol{x})$ , which contains  $\boldsymbol{x}_*$ .

**Remark 5.17.** Suppose that Step 1 yields some vector  $\widehat{\lambda}$ , and let  $\widehat{x}$  be any minimizer of  $f(x) - \widehat{\lambda}^{\top} g(x)$  over all  $x \in \mathbb{R}^{q}$ . If f is strictly convex, then  $\widehat{x} = x_{*}$ . To see this, note that

$$f((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) - \widehat{\boldsymbol{\lambda}}^{\top}\boldsymbol{g}((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) = f((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) - u\widehat{\boldsymbol{\lambda}}^{\top}\boldsymbol{g}(\widehat{\boldsymbol{x}})$$

is a convex function of  $u \in \mathbb{R}$  with value  $L(\widehat{\lambda}) = L(\lambda_*) = f(x_*)$  for  $u \in \{0, 1\}$ , whence

$$f((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) \leq f(\boldsymbol{x}_*) + u\widehat{\boldsymbol{\lambda}}^{\top}\boldsymbol{g}(\widehat{\boldsymbol{x}}) \text{ for } u \in [0,1].$$

On the other hand, by definition of  $L(\widehat{\lambda})$ ,

$$f(\boldsymbol{x}_*) = L(\widehat{\boldsymbol{\lambda}}) \leq f((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) - u\widehat{\boldsymbol{\lambda}}^{\top}\boldsymbol{g}(\widehat{\boldsymbol{x}})$$

for arbitrary  $u \in \mathbb{R}$ , whence

$$f((1-u)\boldsymbol{x}_* + u\widehat{\boldsymbol{x}}) = f(\boldsymbol{x}_*) + u\widehat{\boldsymbol{\lambda}}^{\top}\boldsymbol{g}(\widehat{\boldsymbol{x}}) \text{ for } u \in [0,1].$$

If f is strictly convex, this implies that  $\hat{x} = x_*$ .

**Example 5.18** (Least squares with  $\ell^1$ -penalties). For a given vector  $\boldsymbol{y} \in \mathbb{R}^n$  and a matrix  $\boldsymbol{D} \in \mathbb{R}^{q \times n}$ , consider the function  $H : \mathbb{R}^n \to \mathbb{R}$  with

$$H(\boldsymbol{\beta}) := \|\boldsymbol{y} - \boldsymbol{\beta}\|^2 / 2 + \|\boldsymbol{D}\boldsymbol{\beta}\|_1,$$

where  $\|\cdot\|$  is standard Euclidean norm, and  $\|\boldsymbol{z}\|_1 := \sum_{i=1}^q |z_i|$ . Minimizing H over all vectors  $\boldsymbol{\beta} \in \mathbb{R}^n$  is equivalent to minimizing

$$f(\boldsymbol{\beta}, \boldsymbol{z}) := \| \boldsymbol{y} - \boldsymbol{\beta} \|^2 / 2 + \| \boldsymbol{z} \|_1$$

over all pairs  $(\beta, z) \in \mathbb{R}^n \times \mathbb{R}^q$  under the constraint that  $g(\beta, z) := z - D\beta = 0$ . We may identify  $\mathbb{R}^n \times \mathbb{R}^q$  with  $\mathbb{R}^d$ , where d = n + q, and apply the previous considerations. Note first that for  $\lambda \in \mathbb{R}^q$ ,

$$f(\boldsymbol{\beta}, \boldsymbol{z}) - \boldsymbol{\lambda}^{\top} \boldsymbol{g}(\boldsymbol{\beta}, \boldsymbol{z}) = \|\boldsymbol{y} - \boldsymbol{\beta}\|^2 / 2 + \boldsymbol{\lambda}^{\top} \boldsymbol{D} \boldsymbol{\beta} + \sum_{i=1}^{q} (|z_i| - \lambda_i z_i).$$

For fixed  $\beta$ ,

$$\inf_{\boldsymbol{z}\in\mathbb{R}^{q}}\sum_{i=1}^{q}(|z_{i}|-\lambda_{i}z_{i}) = \begin{cases} 0 & \text{if } \|\boldsymbol{\lambda}\|_{\infty} \leq 1\\ -\infty & \text{if } \|\boldsymbol{\lambda}\|_{\infty} > 1 \end{cases}$$

with  $\|\lambda\|_{\infty} := \max_{i=1,...,q} |\lambda_i|$ . In case of  $\|\lambda\|_{\infty} \le 1$ , the minimum is attained at z if and only if

(5.9) for 
$$i = 1, ..., q$$
,  $z_i \begin{cases} = 0 & \text{if } |\lambda_i| < 1, \\ \ge 0 & \text{if } \lambda_i = 1, \\ \le 0 & \text{if } \lambda_i = -1. \end{cases}$ 

Moreover,

$$\begin{aligned} \|\boldsymbol{y} - \boldsymbol{\beta}\|^2 / 2 + \boldsymbol{\lambda}^\top \boldsymbol{D} \boldsymbol{\beta} &= 2^{-1} \big( \boldsymbol{\beta}^\top \boldsymbol{\beta} - 2 \boldsymbol{\beta}^\top (\boldsymbol{y} - \boldsymbol{D}^\top \boldsymbol{\lambda}) + \|\boldsymbol{y}\|^2 \big) / 2 \\ &\geq \|\boldsymbol{y}\|^2 / 2 - \|\boldsymbol{y} - \boldsymbol{D}^\top \boldsymbol{\lambda}\|^2 / 2 \end{aligned}$$

with equality if and only if

$$\boldsymbol{\beta} = \boldsymbol{y} - \boldsymbol{D}^{\top} \boldsymbol{\lambda}.$$

Hence,

$$egin{aligned} L(oldsymbol{\lambda}) &= \inf_{(oldsymbol{eta},oldsymbol{z})\in\mathbb{R}^n imes\mathbb{R}^q}ig(f(oldsymbol{eta},oldsymbol{z}) - oldsymbol{\lambda}^ opoldsymbol{g}(oldsymbol{eta},oldsymbol{z})ig) \ &= egin{dmatrix} -\infty & ext{if } \|oldsymbol{\lambda}\|_{\infty} > 1, \ \|oldsymbol{y}\|^2 - \|oldsymbol{y} - oldsymbol{D}^ opoldsymbol{\lambda}\|^2/2 & ext{if } \|oldsymbol{\lambda}\|_{\infty} \leq 1, \end{aligned}$$

and in case of  $\|\boldsymbol{\lambda}\|_{\infty} \leq 1$ ,

$$rgmin_{(oldsymbol{eta},oldsymbol{z})\in\mathbb{R}^n imes\mathbb{R}^q}ig(f(oldsymbol{eta},oldsymbol{z})-oldsymbol{\lambda}^ opoldsymbol{g}(oldsymbol{eta},oldsymbol{z})ig)\ =\ ig\{(oldsymbol{y}-oldsymbol{D}^ opoldsymbol{\lambda},oldsymbol{z}):oldsymbol{z}\in\mathcal{Z}(oldsymbol{\lambda})ig\}$$

with  $\mathcal{Z}(\boldsymbol{\lambda})$  denoting the set of all vectors  $\boldsymbol{z} \in \mathbb{R}^q$  satisfying (5.9).

Consequently, to obtain the unique minimizer  $\beta_*$  of H, one may proceed as follows: Find a vector

$$\widehat{oldsymbol{\lambda}} \, \in \, rgmin_{oldsymbol{\lambda} \in \mathbb{R}^q: \|oldsymbol{\lambda}\|_\infty \leq 1} \, \|oldsymbol{y} - oldsymbol{D}^ op oldsymbol{\lambda}\|^2.$$

Then,

 $\widehat{oldsymbol{eta}} \, := \, oldsymbol{y} - oldsymbol{D}^ op \widehat{oldsymbol{\lambda}}$ 

coincides with the unique minimizer  $\beta_*$  of  $H(\beta) = \|\boldsymbol{y} - \beta\|^2/2 + \|\boldsymbol{D}\beta\|_1$  over all  $\beta \in \mathbb{R}^n$ . The reason is that if we would replace  $\hat{\boldsymbol{\lambda}}$  with the Lagrange vector  $\boldsymbol{\lambda}_*$  corresponding to  $\beta_*$ , then  $\boldsymbol{\lambda}_*$  would also minimize  $\|\boldsymbol{y} - \boldsymbol{D}^\top \boldsymbol{\lambda}\|^2$  over all  $\boldsymbol{\lambda} \in \mathbb{R}^q$  satisfying  $\|\boldsymbol{\lambda}\|_{\infty} \leq 1$ , and  $\beta_* = \boldsymbol{y} - \boldsymbol{D}^\top \boldsymbol{\lambda}_*$ . But  $\|\cdot\|^2$  is a strictly convex function on  $\mathbb{R}^n$ , so  $\boldsymbol{D}^\top \hat{\boldsymbol{\lambda}} = \boldsymbol{D}^\top \boldsymbol{\lambda}_*$ , whence  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_*$ .

# 5.2 Special Algorithms

In what follows we consider a continuous and convex function  $f : \mathbb{R}^d \to (-\infty, \infty]$  which is continuously differentiable and strictly convex on  $\mathcal{X} := \operatorname{dom}(f) = \{x \in \mathbb{R}^d : f(x) < \infty\}$ . Furthermore let  $\mathcal{K}$  be a closed convex subset of  $\mathbb{R}^d$  such that

$$\mathcal{X} \cap \mathcal{K} \neq \emptyset,$$

and we assume that  $\{x \in \mathcal{K} : f(x) \leq f(x_o)\}$  is compact for arbitrary  $x_o \in \mathcal{X} \cap \mathcal{K}$ . These conditions imply existence of a unique minimizer

$$oldsymbol{x}_* := rgmin_{oldsymbol{x} \in \mathcal{K}} f(oldsymbol{x})$$

Subsequently we describe algorithms to compute  $x_*$ . Our general results about convex functions show that  $x \in \mathcal{X} \cap \mathcal{K}$  equals  $x_*$  if and only if

(5.10) 
$$\nabla f(\boldsymbol{x})^{\top}(\boldsymbol{y}-\boldsymbol{x}) \geq 0$$
 for arbitrary  $\boldsymbol{y} \in \mathcal{K}$ .

This condition will be encountered repeatedly.

**Exercise 5.19**. Suppose that  $\mathcal{K} = [0, \infty)^d$ . Show that  $x \in \mathcal{X} \cap \mathcal{K}$  equals  $x_*$  if and only if for  $1 \leq j \leq d$ ,

$$\frac{\partial f(\boldsymbol{x})}{\partial x_j} \begin{cases} \geq 0, \\ = 0 & \text{if } x_j > 0. \end{cases}$$

# **5.2.1** Iterative algorithms

As in Section 4.2 we iterate a mapping  $\psi : \mathcal{X} \cap \mathcal{K} \to \mathcal{X} \cap \mathcal{K}$  of type

$$\boldsymbol{\psi}(\boldsymbol{x}) \;=\; \boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x})$$

with a candidate step function  $\Delta : \mathcal{X} \cap \mathcal{K} \to \mathbb{R}^d$  and a step size function  $\lambda : \mathcal{X} \cap \mathcal{K} \to [0, 1]$ . The regularity conditions on  $\Delta$  and  $\lambda$  are essentially identical with the ones in Section 4.2.2:

(i)  $\Delta$  is continuous with  $\Delta(\boldsymbol{x}_*) = \boldsymbol{0}$  while

$$oldsymbol{x}+\Delta(oldsymbol{x})\ \in\ \mathcal{K} \quad ext{and} \quad 
abla f(oldsymbol{x})^{ op}\Delta(oldsymbol{x})\ <\ 0 \quad ext{for}\ oldsymbol{x}\in\mathcal{X}\cap\mathcal{K}\setminus\{oldsymbol{x}_*\}.$$

(ii) With

$$C(\boldsymbol{x}) \ := \ \frac{f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x}))}{\max_{t \in [0,1]} \left[ f(\boldsymbol{x}) - f(\boldsymbol{x} + t\Delta(\boldsymbol{x})) \right]} \quad \text{and} \quad \tilde{C}(\boldsymbol{x}) \ := \ \frac{f(\boldsymbol{x}) - f(\boldsymbol{x} + \lambda(\boldsymbol{x})\Delta(\boldsymbol{x}))}{-\nabla f(\boldsymbol{x})^{\top}\Delta(\boldsymbol{x})}$$

for  $x \in \mathcal{X} \cap \mathcal{K} \setminus \{x_*\}$ , we assume that for any fixed  $y \in \mathcal{X} \cap \mathcal{K} \setminus \{x_*\}$ ,

$$\liminf_{\mathcal{X}\cap\mathcal{K}\ni\boldsymbol{x}\rightarrow\boldsymbol{y}}C(\boldsymbol{x}) > 0 \quad \text{or} \quad \liminf_{\mathcal{X}\cap\mathcal{K}\ni\boldsymbol{x}\rightarrow\boldsymbol{y}}\tilde{C}(\boldsymbol{x}) > 0.$$

With the same arguments as in Section 4.2.2 one can show that iteration of  $\psi$  with arbitrary starting point  $x_0 \in \mathcal{X} \cap \mathcal{K}$  yields a sequence converging to  $\{x_*\}$ .

# A special procedure for $\mathcal{K} = [0, \infty)^d$

In several settings one may parametrize the optimization problem such that the constraints correspond to the convex cone  $\mathcal{K} = [0, \infty)^d$ . Now we derive an explicit example for the candidate

step function  $\Delta$ : As in case of quasi-Newton procedures we approximate f locally by a quadratic function. Precisely let

$$f_{x}(y) := f(x) + \nabla f(x)^{\top} (y - x) + 2^{-1} (y - x)^{\top} A(x) (y - x)$$

for  $\boldsymbol{x} \in \mathcal{X} \cap \mathcal{K}$  with a diagonal matrix

$$\boldsymbol{A}(\boldsymbol{x}) = \operatorname{diag}(a_1(\boldsymbol{x}), \ldots, a_d(\boldsymbol{x}))$$

and continuous diagonal elements  $a_j : \mathcal{X} \cap \mathcal{K} \to (0, \infty)$ . Explicit examples are

$$a_j(oldsymbol{x}) \ := \ 1, \quad a_j(oldsymbol{x}) \ := \ \epsilon(oldsymbol{x}) \ ext{ or } \quad a_j(oldsymbol{x}) \ := \ rac{\partial^2 f(oldsymbol{x})}{\partial x_j^2}$$

In the latter case we assume that these partial derivatives exist and are continuous and strictly positive on  $\mathcal{X} \cap \mathcal{K}$ .

With  $g_j(\boldsymbol{x}) := \nabla f(\boldsymbol{x})_j$  one may write

$$f_{\boldsymbol{x}}(\boldsymbol{y}) = f(\boldsymbol{x}) + \sum_{j=1}^{d} \left( g_j(\boldsymbol{x})(y_j - x_j) + 2^{-1} a_j(\boldsymbol{x})(y_j - x_j)^2 \right)$$
$$= c(\boldsymbol{x}) + 2^{-1} \sum_{j=1}^{d} a_j(\boldsymbol{x}) \left( y_j^2 - 2y_j \left( x_j - \frac{g_j(\boldsymbol{x})}{a_j(\boldsymbol{x})} \right) \right)$$

with c(x) not depending on y. As a function von  $y \in \mathcal{K}$  this is minimal if and only if

$$y_j = \left(x_j - \frac{g_j(\boldsymbol{x})}{a_j(\boldsymbol{x})}\right)^+$$

for  $1 \le j \le d$ . Hence we define

$$\Delta(\boldsymbol{x}) := \underset{\boldsymbol{y} \in \mathcal{K}}{\operatorname{arg\,min}} f_{\boldsymbol{x}}(\boldsymbol{y}) - \boldsymbol{x} = \left( \left( x_j - \frac{g_j(\boldsymbol{x})}{a_j(\boldsymbol{x})} \right)^+ - x_j \right)_{j=1}^d.$$

One can easily verify that this candidate step function  $\Delta$  does fulfill Condition (i).

## An inverse problem from Physics

In a basic experiment in particle physics a random number  $Y \ge 0$  of certain particles is generated, and these particles fly off independently into different directions, the directions being uniformly distributed on the unit sphere of  $\mathbb{R}^3$ . The goal is to estimate the unknown weight function q:  $\mathbb{N}_0 \to [0, \infty)$  with

$$q(y) := \mathbf{I} \mathbf{P}(Y = y),$$

but the total number Y cannot be measured directly. Instead one uses a detector which counts the number  $X \ge 0$  of particles flying off into a direction which is covered by the detector. The relative size  $\rho \in (0,1)$  of these detectable directions is given. So one measurement yields a random variable X such that

$$\mathbb{P}(X = x \,|\, Y = y) = K(x \,|\, y) := \binom{y}{x} \rho^x (1 - \rho)^{y - x},$$

130

and

$$p(x) := \mathbb{P}(X = x) = Kq(x) := \sum_{y=0}^{\infty} K(x \mid y)q(y).$$

Suppose we repeat this basic experiments  $n \gg 1$  times, resulting in independent random counts  $X_1, X_2, \ldots, X_n$  with distribution p. We estimate the weight function p by the empirical weight function  $\hat{p}$  with

$$\widehat{p}(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{[X_i = x]}$$

As shown in Exercise 5.21, the mapping  $q \mapsto p = Kq$  is linear and invertible, but the inverse mapping is ill-conditioned in the sense that small deviations of  $\hat{p}$  from p may lead to rather large errors in the reconstruction of q.

An alternative approach is to minimize the functional

$$f(q) \ := \ -\sum_{x=0}^{\infty} \widehat{p}(x) \log Kq(x) + \sum_{y=0}^{\infty} q(y)$$

over all weight functions  $q : \mathbb{N}_0 \to [0, \infty)$ . Indeed, one can show that in case of  $\hat{p} \equiv p = Kq$ , the unique minimizer of f is given by q. In general, any minimizer q of f satisfies automatically the constraint  $\sum_y q(y) = 1$ , because for any weight function q with  $f(q) < \infty$  and arbitrary numbers t > 0,

$$\frac{d}{dt}f(tq) = \sum_{y=0}^{\infty} q(y) - t^{-1}$$

Note that for integers  $0 \le x \le y$ ,

$$\frac{K(x\,|\,y+1)}{K(x\,|\,y)} \;=\; \frac{1-\rho}{1-x/(y+1)} \;<\; 1 \quad \text{if and only if} \quad y \;>\; x/\rho-1.$$

Hence, if we define  $\hat{x} := \max(X_1, \ldots, X_n)$  and

$$\widehat{y} := \begin{cases} 0 & \text{if } \widehat{x} = 0\\ \lfloor \widehat{x}/\rho \rfloor - 1 & \text{if } \widehat{x} > 0 \end{cases}$$

for any weight function q and  $x \leq \hat{x}$ , the value Kq(x) gets strictly larger if we replace

$$\big(q(\widehat{y}),q(\widehat{y}+1),q(\widehat{y}+2),q(\widehat{y}+3),\ldots\big) \quad \text{with} \quad \Big(\sum_{y \geq \widehat{y}} q(y),0,0,0,\ldots\Big),$$

unless q(y) = 0 for all  $y > \hat{y}$ . Consequently, we have to solve a finite-dimensional minimization problem:

With the dimensions  $c := \hat{x} + 1$  and  $d := \hat{y} + 1$ , the vector  $\hat{p} := (\hat{p}(i-1))_{i=1}^{c}$  and the matrix  $K := (K(i-1|j-1))_{1 \le i \le c, 1 \le j \le d}$  we want to minimize

$$f(\boldsymbol{q}) := -\sum_{i=1}^{c} \widehat{p}_i \log(\boldsymbol{K}\boldsymbol{q})_i + \sum_{j=1}^{d} q_j$$

over all vectors  $q \in [0, \infty)^d$ , where q corresponds to  $(q(j-1))_{j=1}^d$ . Note that with the convention  $\log(r) := -\infty$  for  $r \leq 0$ , the function  $f : \mathbb{R}^d \to (-\infty, \infty]$  is continuous with  $f(q) < \infty$  if and only if

$$(\mathbf{Kq})_i > 0$$
 whenever  $\widehat{p}_i > 0$ .

Moreover, for  $\boldsymbol{q} = t\boldsymbol{u}$  with  $t \ge 0$  and  $\boldsymbol{u} \in \Sigma_d := \{ \boldsymbol{v} \in [0,1]^d : \sum_{j=1}^d v_j = 1 \}$ ,

$$f(\boldsymbol{q}) = -\log t + t - 1 + f(\boldsymbol{u}) \geq -\log t + t - 1 + \min_{\boldsymbol{v} \in \Sigma_d} f(\boldsymbol{v}) \to \infty \quad \text{as } t \to \{0, \infty\}.$$

Hence the set  $\{x \in [0,\infty)^d : f(x) \le f(q)\}$  is compact for any  $q \in [0,\infty)^d \cap \operatorname{dom}(f)$ . Finally, f is twice continuously differentiable on  $\operatorname{dom}(f)$  with

$$g_j(\boldsymbol{q}) := \frac{\partial}{\partial q_j} f(\boldsymbol{q}) = 1 - \sum_{i=1}^c \frac{\widehat{p}_i K_{ij}}{(\boldsymbol{K}\boldsymbol{q})_i},$$
$$a_j(\boldsymbol{q}) := \frac{\partial^2}{\partial q_j^2} f(\boldsymbol{q}) = \sum_{i=1}^c \frac{\widehat{p}_i K_{ij}^2}{(\boldsymbol{K}\boldsymbol{q})_i^2}.$$

Hence we may apply the general optimization algorithm described before.

**Remark 5.20** (EM algorithm). For the particular problem here, many statisticians would use a so-called expectation-minimization algorithm. Suppose we would observe independent pairs  $(X_1, Y_1), \ldots, (X_n, Y_n)$  with  $Y_i$  having distribution q and  $\mathbb{P}(X_i = x | Y_i = y) = K(x | y)$ . Then one could estimate q by the unobserved empirical weight function  $\check{q}$  with

$$\check{q}(y) := \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{[Y_i = y]}$$

Indeed,  $\check{q}$  is the unique minimizer of

$$\check{f}(q) \; := \; -\sum_{y=0}^{\infty}\check{q}(y)\log q(y) + \sum_{y=0}^{\infty}q(y)$$

Since we observe only  $X_1, \ldots, X_n$ , the idea is to replace a current estimator  $q_k$  for the true q by

$$q_{k+1}(y) := \mathbb{E}_{q_k} \left( \check{q}(y) \mid X_1, \dots, X_n \right)$$
$$= \sum_{x=0}^{\widehat{x}} \widehat{p}(x) \mathbb{P}_{q_k} (Y_1 = y \mid X_1 = x)$$
$$= \sum_{x=0}^{\widehat{x}} \frac{\widehat{p}(x) K(x \mid y) q_k(y)}{\sum_{z \ge 0} K(x \mid z) q_k(z)}.$$

Here  $\mathbb{E}_{q_k}$  and  $\mathbb{P}_{q_k}$  denote expectations and probabilities in case of  $q \equiv q_k$ . Although this approach is rather intuitive and easy to implement, the resulting algorithm converges slower than the algorithm we derived before.

**Exercise 5.21**. Show that for  $s \in \mathbb{R}$  and  $k \in \mathbb{N}_0$ ,

$$\mathbb{E}([X]_k s^{X-k}) = \rho^k \mathbb{E}([Y]_k (1-\rho+\rho s)^{Y-k}),$$

provided that both sides are well-defined in  $\mathbb{R}$ . Here  $[z]_0 := 1$  and  $[z]_k := \prod_{i=0}^{k-1} (z-i)$  for  $z \in \mathbb{R}$  and  $k \in \mathbb{N}$ .

Deduce from the previous identities that

$$\mathbb{E}(Y) = \rho^{-1} \mathbb{E}(X)$$
 and  $Var(Y) = \rho^{-2} Var(X) - \rho^{-2}(1-\rho) \mathbb{E}(X).$ 

Furthermore, show that for arbitrary  $k \in \mathbb{N}_0$ ,

$$q(k) = \frac{1}{\rho^k k!} \mathbb{E}([X]_k (1 - 1/\rho)^{X-k}),$$

provided that  $\mathbb{E}(s_o^Y) < \infty$  for some  $s_o > 2(1 - \rho)$ .

# 5.2.2 Active set methods

In this section we concentrate on a polyhedral cone

$$\mathcal{K} = \left\{ \boldsymbol{x} \in \mathbb{R}^d : \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle \le 0 \text{ for } i = 1, \dots, q \right\}$$

with  $q \leq d$  linearly independent vectors  $v_1, \ldots, v_q \in \mathbb{R}^d$ . For a vector  $x \in \mathbb{R}^d$  let

$$J(\boldsymbol{x}) := \{ i \in \{1, \ldots, q\} : \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle = 0 \}.$$

If  $x \in \mathcal{K}$ , then J(x) specifies the "set of active constraints (active set) at x", that is, the set of inequality constraints which are even equalities for the point x. For an arbitrary index set  $J \subset \{1, \ldots, q\}$  let

$$\mathcal{X}_J := \mathcal{X} \cap \{ \boldsymbol{x} \in \mathbb{R}^d : \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle = 0 \text{ for all } i \in J \},$$

i.e. the intersection of the open set  $\mathcal{X} = \text{dom}(f)$  with a (d - #J)-dimensional linear subspace of  $\mathbb{R}^d$ . That is,  $\boldsymbol{x} \in \mathcal{X}_J$  if and only if  $J(\boldsymbol{x}) \supset J$ . Moreover,  $\boldsymbol{x} \in \mathcal{X}_J$  belongs to  $\mathcal{K}$  if and only if

$$\langle \boldsymbol{x}, \boldsymbol{v}_i \rangle \leq 0 \quad \text{for all } i \in \{1, \dots, q\} \setminus J.$$

The idea of active set methods is to determine for an index set  $J \subset \{1, \ldots, q\}$  the point

$$oldsymbol{x}_J := rgmin_{\mathcal{X}_J} f$$

(if possible) and to modify the set J finitely many times until  $x_J = \{x_*\}$  eventually. Here we assume that minimizing f over such a set  $\mathcal{X}_J$  is sufficiently easy. In case of a quadratic function f this is certainly the case, see Section 5.1.2.

An obvious question is how to infer for a given index set J that the vector  $x_J$  is equal to  $x_*$ . To this end we chose a "dual basis"  $b_1, b_2, \ldots, b_d$  of  $\mathbb{R}^d$  in the sense that for  $i \in \{1, \ldots, d\}$  and  $j \in \{1, \ldots, q\}$ ,

$$\langle \boldsymbol{b}_i, \boldsymbol{v}_j \rangle = 1_{[i=j]}.$$

(Such a dual basis always exists. For instance, one may augment the given vectors  $v_1, \ldots, v_q$  to a basis  $v_1, \ldots, v_d$  of  $\mathbb{R}^d$ , and then the rows of the matrix  $[b_1, \ldots, b_d]^\top := [v_1, \ldots, v_d]^{-1}$  have the desired property.) With these basis vectors  $b_i$  we may write

$$\mathcal{K} = \Big\{ \sum_{i=1}^{d} \lambda_i \boldsymbol{b}_i : \boldsymbol{\lambda} \in \mathbb{R}^d, \lambda_i \leq 0 \text{ for } i = 1, \dots, q \Big\},$$
  
$$\mathcal{X}_J = \mathcal{X} \cap \Big\{ \sum_{i=1}^{d} \lambda_i \boldsymbol{b}_i : \boldsymbol{\lambda} \in \mathbb{R}^d, \lambda_i = 0 \text{ for all } i \in J \Big\}.$$

**Lemma 5.22** (Characterizing  $x_J$  and  $x_*$ ). Let  $J \subset \{1, \ldots, q\}$  and  $x \in \mathcal{X}_J$ .

(a) The vector x equals  $x_J$  if and only if

(5.11) 
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_i \rangle = 0 \text{ for all } i \in \{1, \dots, d\} \setminus J.$$

(b) If  $x \in \mathcal{K}$  and satisfies (5.11), then it equals  $x_*$  if and only if

(5.12) 
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_j \rangle \leq 0 \quad \text{for all } j \in J.$$

**Proof of Lemma 5.22.** Part (a) is an easy consequence of convexity of f plus the fact that  $\mathcal{X}_J$  is the intersection of the open set  $\mathcal{X}$  with the linear subspace span  $\{\mathbf{b}_i : i \in \{1, \ldots, d\} \setminus J\} \subset \mathbb{R}^d$ .

Concerning part (b), suppose first that  $x = x_*$ . For arbitrary  $j \in J$ , the vector  $y = x - b_j$  lies in  $\mathcal{K}$ , so by (5.10),

$$0 \leq \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle = - \langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_j \rangle.$$

Hence x has to fulfill (5.12).

On the other hand, suppose that  $\boldsymbol{x}$  lies within  $\mathcal{K}$  and satisfies both (5.11) and (5.12). Any vector  $\boldsymbol{y} \in \mathcal{K}$  may be written as  $\boldsymbol{x} + \sum_{i=1}^{d} \lambda_i \boldsymbol{b}_i$ , where  $\lambda_i \leq 0$  in case of  $\langle \boldsymbol{x}, \boldsymbol{b}_i \rangle = 0$ . In particular,  $\lambda_i \leq 0$  for all  $i \in J$ . Thus,

$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle = \sum_{i \in J} \underbrace{\lambda_i}_{\leq 0} \underbrace{\langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_i \rangle}_{\leq 0} + \sum_{i \in \{1, \dots, d\} \setminus J} \lambda_i \underbrace{\langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_i \rangle}_{= 0} \geq 0.$$

Consequently, x satisfies (5.10) and hence is equal to  $x_*$ .

# An active set algorithm for "quadratic programming"

We consider the special case of a quadratic function f with positive definite Hessian matrix. The active set algorithm alternates finitely many times between two basic procedures. Both precedures return a pair (x, F) consisting of a candidate  $x \in \mathcal{K}$  for  $x_*$  and the value F = f(x).

An important notion of the algorithm is *local optimality* of a candidate  $x \in \mathcal{K}$  for  $x_*$ . This means that

$$\boldsymbol{x} = \boldsymbol{x}_{J(\boldsymbol{x})}$$

That is,  $\boldsymbol{x}$  lies in  $\mathcal{K}$  and minimizes f over the linear subspace  $\mathcal{X}_{J(\boldsymbol{x})}$  of all points  $\tilde{\boldsymbol{x}} \in \mathbb{R}^d$  such that  $J(\tilde{\boldsymbol{x}}) \supset J(\boldsymbol{x})$ .

**Starting point.** A natural starting point is  $x_{\{1,...,\}}$ , because  $\mathcal{X}_{\{1,...,q\}} \subset \mathcal{K}$ . This point is clearly locally optimal.

**Basic procedure 1 (checking optimality; deactivating constraints).** Suppose that  $x \in \mathcal{K}$  is locally optimal. Then we check Condition (5.12) in Lemma 5.22: If  $\langle \nabla f(x), \mathbf{b}_i \rangle \leq 0$  for all  $i \in J(x)$ , we know that  $x = x_*$  and stop the algorithm. Otherwise we replace x with some vector  $x_{\text{new}} \in \mathcal{K}$  such that  $f(x_{\text{new}}) < f(x)$ .

Here is a particular variant for  $\boldsymbol{x}_{new}$ : For some nonvoid set  $\tilde{J}(\boldsymbol{x}) \subset J(\boldsymbol{x})$  such that  $\langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_i \rangle > 0$  for all  $i \in \tilde{J}(\boldsymbol{x})$ , we compute

$$oldsymbol{x}_{ ext{new}} \ = \ oldsymbol{x} - \sum_{i \in ilde{J}(oldsymbol{x})} \langle 
abla f(oldsymbol{x}), oldsymbol{b}_i 
angle oldsymbol{b}_i.$$

This vector is certainly in  $\mathcal{K}$ . Then we check whether  $f(\boldsymbol{x}_{\text{new}}) < f(\boldsymbol{x})$ . As long as this is not the case, we replace  $\boldsymbol{x}_{\text{new}}$  with  $2^{-1}(\boldsymbol{x} + \boldsymbol{x}_{\text{new}}) \in \mathcal{K}$ . After finitely many steps,  $f(\boldsymbol{x}_{\text{new}}) < f(\boldsymbol{x})$ , because

$$\lim_{t \to 0+} t^{-1} \big( f((1-t)\boldsymbol{x} + t\boldsymbol{x}_{\text{new}}) - f(\boldsymbol{x}) \big) = \langle \nabla f(\boldsymbol{x}), \boldsymbol{x}_{\text{new}} - \boldsymbol{x} \rangle = -\sum_{i \in \tilde{J}(\boldsymbol{x})} \langle \nabla f(\boldsymbol{x}), \boldsymbol{b}_i \rangle^2$$

is strictly negative. The fact that a finite number of updates  $x_{\text{new}} \leftarrow (x + x_{\text{new}})/2$  is sufficient is true in theory. In practice, if a moderate number of updates does not yield a better candidate  $x_{\text{new}}$ , we conclude that x is equal to  $x_*$  up to numerical errors.

**Basic procedure 2 (local search).** Starting from an arbitrary candidate  $x \in \mathcal{K}$  for  $x_*$ , we want to replace it with a point  $x_{new} \in \mathcal{K}$  which is locally optimal and satisfies

$$(5.13) J(\boldsymbol{x}_{\text{new}}) \supset J(\boldsymbol{x}),$$

$$(5.14) f(\boldsymbol{x}_{\text{new}}) \leq f(\boldsymbol{x}).$$

To this end, we determine J := J(x) and compute  $x_{new} := x_J$ . In case of  $f(x_{new}) = f(x)$  we know that  $x_{new} = x$  is locally optimal, and requirements (5.13), (5.14) are obviously satisfied.

If  $f(\boldsymbol{x}_{\text{new}}) < f(\boldsymbol{x})$ , we know that  $f((1-t)\boldsymbol{x} + t\boldsymbol{x}_{\text{new}})$  is strictly decreasing in  $t \in [0, 1]$ . Note that by definition of  $\boldsymbol{x}_{\text{new}}$ ,

$$\langle \boldsymbol{x}, \boldsymbol{v}_i \rangle = 0 = \langle \boldsymbol{x}_{\text{new}}, \boldsymbol{v}_i \rangle$$
 for all  $i \in J$ ,

while the definition of J implies that

$$\langle \boldsymbol{x}, \boldsymbol{v}_i \rangle < 0$$
 whenever  $i \in \{1, \ldots, q\} \setminus J$ .

Thus,  $oldsymbol{x}_{ ext{new}} \in \mathcal{K}$  if and only if

$$V(J, \boldsymbol{x}_{ ext{new}}) := \left\{ i \in \{1, \dots, q\} \setminus J : \langle \boldsymbol{x}_{ ext{new}}, \boldsymbol{v}_i \rangle > 0 \right\}$$

```
(x, F) \leftarrow \text{BasicProcedure2}(x)
F \leftarrow f(\boldsymbol{x})
                                                                                             (**)
J \leftarrow J(\boldsymbol{x})
oldsymbol{x}_{	ext{new}} \leftarrow oldsymbol{x}_{.I}
F_{\text{new}} \leftarrow f(\boldsymbol{x}_{\text{new}})
while F_{\text{new}} < F do
                                                                                              (**)
               V \leftarrow V(J, \boldsymbol{x}_{\text{new}})
               if V = \emptyset then
                              x \leftarrow x_{	ext{new}}
                              F \leftarrow F_{\text{new}}
               else
                              t_o \leftarrow t_o(J, \boldsymbol{x}_{\text{new}})
                              \boldsymbol{x}_{\text{new}} \leftarrow (1 - t_o)\boldsymbol{x} + t_o \boldsymbol{x}_{\text{new}}
                              F_{\text{new}} \leftarrow f(\boldsymbol{x}_{\text{new}})
                                                                                                (*)
                              if F_{\text{new}} < F then
                                             m{x} \leftarrow m{x}_{	ext{new}}
                                              F \leftarrow F_{\text{new}}
                                                                                               (**)
                                             J \leftarrow J(\boldsymbol{x})
                                             x_{	ext{new}} \leftarrow x_{I}
                                             F_{\text{new}} \leftarrow f(\boldsymbol{x}_{\text{new}})
                              end if
               end if
end while.
```

Table 5.1: Basic procedure 2 for quadratic f.

is empty. In that case, it follows from  $J(\boldsymbol{x}_{\text{new}}) \supset J$  that  $\boldsymbol{x}_{\text{new}}$  is locally optimal, satisfies (5.13) and satisfies (5.14) with strict inequality.

If  $V(J, \boldsymbol{x}_{\text{new}}) \neq \emptyset$ , the largest value  $t \in [0, 1]$  such that  $(1 - t)\boldsymbol{x} + t\boldsymbol{x}_{\text{new}} \in \mathcal{K}$  is given by

$$t_o(J, \boldsymbol{x}_{\mathrm{new}}) = \min_{i \in V(J, \boldsymbol{x}_{\mathrm{new}})} rac{-\langle \boldsymbol{x}, \boldsymbol{v}_i 
angle}{\langle \boldsymbol{x}_{\mathrm{new}}, \boldsymbol{v}_i 
angle - \langle \boldsymbol{x}, \boldsymbol{v}_i 
angle} \in (0, 1).$$

If we set  $t_o := t_o(J, \boldsymbol{x}_{new})$  and replace  $\boldsymbol{x}_{new}$  with  $(1 - t_o)\boldsymbol{x} + t_o\boldsymbol{x}_{new}$ , then the new point  $\boldsymbol{x}_{new}$  belongs to  $\mathcal{K}$ , satisfies (5.13) with strict inclusion and (5.14) with strict inequality. But it may still fail to be locally optimal. Thus we replace  $\boldsymbol{x}$  with  $\boldsymbol{x}_{new}$  and repeat the previous steps.

After a finite number of iterations, the point  $\boldsymbol{x}_{new}$  will be locally optimal and satisfy the requirements (5.13) and (5.14). Indeed, whenever we replace  $\boldsymbol{x}$  with a new vector, the value  $f(\boldsymbol{x})$  decreases strictly, the set  $J(\boldsymbol{x})$  increases strictly, so eventually the point  $\boldsymbol{x}_{new} = \boldsymbol{x}_{J(\boldsymbol{x})}$  will have to satisfy  $V(\boldsymbol{x}, \boldsymbol{x}_{new}) = \emptyset$ .

Table 5.1 provides pseudo-code for the whole procedure. Checking the inequality  $F_{\text{new}} < F$  in line (\*) is superflous in theory, but in practice it could be violated due to rounding errors. In that case we conclude that  $\boldsymbol{x}$  is already locally optimal up to rounding errors. Whenever one determines the set  $J = J(\boldsymbol{x})$  or the set  $V(J, \boldsymbol{x}_{\text{new}})$  in lines (\*\*), one should be wary or numerical errors and replace conditions such as  $\pm \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle > 0$  with  $\pm \langle \boldsymbol{x}, \boldsymbol{v}_i \rangle > \tau$  with a very small parameter  $\tau > 0$ . Without such precautions, there is a risk of creating an endless loop.



Table 5.2: Active set algorithm for quadratic f.

**Complete algorithm.** Table 5.2 describes the whole algorithm as described here.

# **Example: Concave regression**

Suppose that for given data vectors  $X, Y \in \mathbb{R}^n$ , we are looking for a concave function  $\hat{\mu} : \mathbb{R} \to \mathbb{R}$  such that

$$\sum_{\ell=1}^{n} (Y_{\ell} - \widehat{\mu}(X_{\ell}))^2$$

is minimized. The minimizer  $\hat{\mu}$  is not unique, but it is unique on the set  $\{X_1, \ldots, X_n\}$ . Precisely, let

 $t_1 < t_2 < \dots < t_d$ 

be the different elements of  $\{X_1, \ldots, X_n\}$ , and let

$$w_i := \#\{\ell \le n : X_\ell = x_i\}.$$

Then for any function  $\widehat{\mu} : \mathbb{R} \to \mathbb{R}$ ,

$$\sum_{\ell=1}^{n} (Y_{\ell} - \widehat{\mu}(X_{\ell}))^2 = S_0^2 + \sum_{i=1}^{d} w_i (y_i - \widehat{\mu}(t_i))^2,$$

where

$$S_0^2 := \sum_{i=1}^d \sum_{\ell: X_\ell = x_i} (Y_\ell - y_i)^2 = \|\mathbf{Y}\|^2 - \sum_{i=1}^d w_i y_i^2$$

Thus, we are looking for a vector  $m \in \mathbb{R}^d$ , interpreted as  $(\widehat{\mu}(t_i))_{i=1}^d$ , with minimal value of

$$f(\boldsymbol{m}) := \sum_{i=1}^{d} w_i (y_i - m_i)^2$$

under the constraints that

(5.15) 
$$\frac{m_k - m_{k-1}}{t_k - t_{k-1}} \ge \frac{m_{k+1} - m_k}{t_{k+1} - t_k} \quad \text{for } 1 < k < d.$$

A possible extension of m to a concave function  $\hat{\mu}$  on  $\mathbb{R}$  is to require that  $\hat{\mu}$  is affine on each interval  $[t_i, t_{i+1}], 1 \leq i < d$ , and on  $(-\infty, t_2]$  as well as on  $[t_{d-1}, \infty)$ . (It is tacitly assumed that  $d \geq 3$ .)

In this particular example it is convenient to label the constraint vectors for  $\mathcal{K}$  with indices  $k \in \{2, \ldots, d-1\}$ . Precisely, let

$$\begin{aligned} \boldsymbol{v}_k &:= \Big(\frac{1_{[i=k+1]} - 1_{[i=k]}}{t_{k+1} - t_k} - \frac{1_{[i=k]} - 1_{[i=k-1]}}{t_k - t_{k-1}}\Big)_{i=1}^d \\ &= \Big(\frac{(t_{k+1} - t_k)1_{[i=k-1]} - (t_{k+1} - t_{k-1})1_{[i=k]} + (t_k - t_{k-1})1_{[i=k+1]}}{(t_{k+1} - t_k)(t_k - t_{k-1})}\Big)_{i=1}^d \end{aligned}$$

for  $2 \le k \le d - 1$ . Then (5.15) is equivalent to

$$\langle \boldsymbol{m}, \boldsymbol{v}_k \rangle \leq 0 \quad \text{for } 2 \leq k \leq d-1$$

A suitable dual basis is given, for instance, by the vectors

$$b_1 := (1)_{i=1}^d, b_k := ((t_k - t_i)^+)_{i=1}^d, \quad 2 \le k \le d.$$

For a set  $J \subset \{2, \ldots, d-1\}$ , the set

$$\mathcal{X}_J = \left\{ \boldsymbol{m} \in \mathbb{R}^d : \langle \boldsymbol{m}, \boldsymbol{v}_k \rangle = 0 \text{ for all } k \in J \right\}$$

corresponds to the set of all continuous and piecewise linear functions  $\mu : [t_1, t_d] \to \mathbb{R}$  with potential changes of slope only at the points  $t_i, i \in \{1, \ldots, d\} \setminus J$ . These points will be referred to as *knots*.

Concerning the set  $\tilde{J}(\boldsymbol{m})$  for basic procedure 1, in our particular implementation, we go through all pairs of consecutive knots  $x_k$  and  $x_\ell$  of  $\boldsymbol{m}$  such that  $k + 1 < \ell$ . Then we choose an index

$$i(k,\ell) \in rgmax_{i\in\{k+1,\ldots,\ell-1\}} \langle \nabla f(\boldsymbol{m}), \boldsymbol{b}_i \rangle.$$

If  $\langle \nabla f(\boldsymbol{m}), \boldsymbol{b}_{i(k,\ell)} \rangle$  is larger than 0.5 times  $\max_{i=1,\dots,d} \langle \nabla f(\boldsymbol{m}), \boldsymbol{b}_i \rangle$ , then  $i(k,\ell) \in \tilde{J}(\boldsymbol{m})$ .

One important fact is that we don't have to generate and store the vectors  $v_k$  and  $b_i$ . Indeed, fitting a continuous and piecewise linear function with specified knots (a linear spline) is a standard problem from linear regression. Moreover, for

$$\boldsymbol{g} := \nabla f(\boldsymbol{m}) = -2 \big( w_i (y_i - m_i) \big)_{i=1}^d$$

and  $k \in \{2, \ldots, d-1\}$  one may write

$$\langle \boldsymbol{g}, \boldsymbol{b}_k \rangle = \sum_{i=1}^{a} g_i (t_k - t_i)^+$$
  
=  $\sum_{i=1}^{k-1} g_i \sum_{j=i}^{k-1} (t_{j+1} - t_j)$   
=  $\sum_{j=1}^{k-1} (t_{j+1} - t_j) \sum_{i=1}^{j} g_i.$ 

Hence one may compute all directional derivatives  $\langle g, b_k \rangle$ ,  $2 \leq k < d$ , by a simple summation scheme.

Figures 5.2 and 5.3 illustrate this procedure with a data example involving n = 1000 raw observations  $(X_{\ell}, Y_{\ell})$ , leading to d = 101 triples  $(t_i, y_i, w_i)$ . The raw observations  $(X_{\ell}, Y_{\ell})$  and pairs  $(t_i, y_i)$  are shown in the upper panel of Figure 5.2. The lower panel shows the final fit. Figure 5.3 shows intermediate stages of the active set algorithm. The current fit is shown as a blue line with the knots (changes of slope) highlighted. On sees the current fit after one, three, five and seven iterations (i.e. basic procedure 1, followed by basic procedure 2). That is, the current fit is locally optimal. The vertical lines (magenta) indicate positions at which the currently active constraint will be deactivated, leading to a new knot.

# Active set algorithms for general f

For the general setting, suppose that we work with quasi-Newton steps for the miminization of f over  $\mathcal{X}_J$ . That is, in addition to f and  $\nabla f$  there is a continuous mapping  $A : \mathcal{X} \to \mathbb{R}^{d \times d}_{\text{sym},+}$ , and for a candidate  $x_o \in \mathcal{K} \cap \mathcal{X}$  for  $x_*$ , the function f(x) is approximated temporarily by the quadratic function

$$\widetilde{f}(\boldsymbol{x} \mid \boldsymbol{x}_o) := f(\boldsymbol{x}_o) + \nabla f(\boldsymbol{x}_o)^\top (\boldsymbol{x} - \boldsymbol{x}_o) + 2^{-1} (\boldsymbol{x} - \boldsymbol{x}_o)^\top A(\boldsymbol{x}_o) (\boldsymbol{x} - \boldsymbol{x}_o).$$

Note that  $\boldsymbol{x}_o = \boldsymbol{x}_*$  if and only if  $\boldsymbol{x}_o$  minimizes  $\tilde{f}(\cdot | \boldsymbol{x}_o)$  over  $\mathcal{K}$ .

One strategy, used by numerous authors, is to apply an active set algorithm for quadratic functions to this surrogate function  $\tilde{f}(\cdot | \mathbf{x}_o)$  to obtain a new candidate  $\psi_o(\mathbf{x}_o) \in \mathcal{K}$  for  $\mathbf{x}_*$ . Then one applies some stepsize correction and determines a point  $\psi(\mathbf{x}_o) = (1 - \lambda(\mathbf{x}_o))\mathbf{x}_o + \lambda(\mathbf{x}_o)\psi_o(\mathbf{x}_o)$  to make sure that this point does belong to  $\mathcal{K} \cap \mathcal{X}$  and yields a strictly smaller value of f, unless  $\mathbf{x}_o$  was already the minimizer  $\mathbf{x}_*$ .

# 5.2.3 Isotonic least squares regression

The present and subsequent subsection is about the special cone

$$\mathcal{K}_{\uparrow} := \left\{ oldsymbol{m} \in \mathbb{R}^d : m_1 \leq m_2 \leq \cdots \leq m_d 
ight\}.$$

At first we'll discuss a particular optimization problem on this set  $\mathcal{K}_{\uparrow}$  and describe geometrical properties of the solution. In the next subsection we consider a rather general optimization problem

and provide an exact algorithm for its solution which may be viewed as a special active set method. All this material is based on ideas of Robertson and Waltman (1968) and Robertson et al. (1988). For a fixed data vector  $\boldsymbol{y} \in \mathbb{R}^d$  and a weight vector  $\boldsymbol{w} \in (0, \infty)^d$  we consider the function

$$f(\boldsymbol{m}) := \sum_{i=1}^d w_i (y_i - m_i)^2$$

Obviously this quadratic function is strictly convex and coercive on  $\mathbb{R}^d$ , i.e.  $f(\boldsymbol{m}) \to \infty$  as  $\|\boldsymbol{m}\| \to \infty$ . Consequently,  $\boldsymbol{x}_* = \arg\min_{\boldsymbol{m}\in\mathcal{K}_{\uparrow}} f(\boldsymbol{m})$  exists. Here

$$\nabla f(\boldsymbol{m}) = 2 \left( w_i (m_i - y_i) \right)_{i=1}^d,$$

so it follows from (5.10) that the minimizer of f over  $\mathcal{K}_{\uparrow}$  is the unique point  $m \in \mathcal{K}_{\uparrow}$  such that

(5.16) 
$$\sum_{i=1}^{d} w_i (m_i - y_i) \Delta_i \geq 0 \quad \text{whenever } \Delta \in \mathbb{R}^d \text{ such that } \boldsymbol{m} + t \Delta \in \mathcal{K}_{\uparrow} \text{ for some } t > 0.$$

Plugging in suitable test vectors  $\Delta$  in the previous inequality yieds the following characterization of  $x_*$ .

**Theorem 5.23.** A vector  $m \in \mathcal{K}_{\uparrow}$  minimizes the function f above over  $\mathcal{K}_{\uparrow}$  if and only if for arbitrary  $j \in \{1, 2, ..., d\}$ ,

$$\sum_{i=1}^{j} w_i m_i \leq \sum_{i=1}^{j} w_i y_i$$

with equality if j = d or  $m_j < m_{j+1}$ .

**Remark 5.24** (Convex minorants). Before proving the latter theorem, let us interpret it geomatrically. An arbitrary vector  $v \in \mathbb{R}^d$  may be identified with the picewise linear path  $S(v) \subset \mathbb{R} \times \mathbb{R}$ , connecting the d + 1 points (0, 0) and

$$\left(\sum_{i=1}^{j} w_i, \sum_{i=1}^{j} w_i v_i\right), \quad 1 \le j \le d$$

A vector  $\boldsymbol{m}$  belongs to  $\mathcal{K}_{\uparrow}$  if and only if its path  $\mathcal{S}(\boldsymbol{m})$  is the graph of a convex function on the interval  $[0, \sum_{i=1}^{f} w_i]$ . According to Theorem 5.23, this path has to be below the path  $\mathcal{S}(\boldsymbol{y})$  with the same starting point and end point. Moreover,  $m_j < m_{j+1}$  is equivalent to saying that the path  $\mathcal{S}(\boldsymbol{m})$  changes slope at the point  $(\sum_{i=1}^{j} w_i, \sum_{i=1}^{j} w_i m_i)$ . Consequently, at any such point the path  $\mathcal{S}(\boldsymbol{m})$  has to touch the path  $\mathcal{S}(\boldsymbol{y})$ . Now one can verify easily that this amounts to the following description: The path  $\mathcal{S}(\boldsymbol{m})$  corresponds to the largest (pointwise) convex function on the interval  $[0, \sum_{i=1}^{f} w_i]$  which is less than or equal to the function corresponding to the path  $\mathcal{S}(\boldsymbol{y})$ .

This fact is illustrated in Figure 5.4. In the upper panel one sees a scatter plot of d = 40 pairs  $(t_i, y_i) = (i, y_i)$  as well as a line plot of the optimal vector  $\mathbf{m} \in \mathcal{K}_{\uparrow}$ , where we chose  $w_i = 1$ . In the lower panel one sees the corresponding paths  $\mathcal{S}(\mathbf{y})$  and  $\mathcal{S}(\mathbf{m})$ .

**Proof of Theorem 5.23.** Let m be the minimizer of f on  $\mathcal{K}_{\uparrow}$ . Clearly the vector  $\Delta := \pm (1)_{i=1}^d$  satisfies the condition in (5.16). Hence  $\pm \sum_{i=1}^d w_i(m_i - y_i) \ge 0$  which is equivalent to

$$\sum_{i=1}^d w_i m_i = \sum_{i=1}^d w_i y_i.$$

For  $1 \le j < d$  we consider the vector  $\Delta := (-1_{[i \le j]})_{i=1}^d$ . It satisfies the condition in (5.16), too, whence  $-\sum_{i \le j} w_i(m_i - y_i) \ge 0$  or, equivalently,

$$\sum_{i=1}^j w_i m_i \leq \sum_{i=1}^j w_i y_i.$$

In case of  $m_j < m_{j+1}$  the vector  $\Delta := (1_{[i \leq j]})_{i=1}^d$  fulfils the same requirement, whence

$$\sum_{i=1}^j w_i m_i = \sum_{i=1}^j w_i y_i.$$

For the reverse direction, let  $m \in \mathcal{K}_{\uparrow}$  such that  $\sum_{i=1}^{j} w_i m_i \leq \sum_{i=1}^{j} w_i y_i$  for all  $j \in \{1, \ldots, d\}$  with equality in case of j = d or  $m_j < m_{j+1}$ . Now let  $\Delta \in \mathbb{R}^d$  be a vector such that  $m + t\Delta \in \mathcal{K}_{\uparrow}$  for some t > 0. This is equivalent to saying that

$$(5.17) \qquad \qquad \Delta_j \le \Delta_{j+1} \quad \text{if } 1 \le j < d \text{ and } m_j = m_{j+1}$$

Since  $\sum_{i=1}^{d} w_i(m_i - y_i) = 0$ , we may write

$$\sum_{i=1}^{d} w_i(m_i - y_i)\Delta_i = \sum_{i=1}^{d} w_i(m_i - y_i)(\Delta_i - \Delta_d)$$
  
= 
$$\sum_{i=1}^{d-1} w_i(m_i - y_i) \sum_{j=i}^{d-1} (\Delta_j - \Delta_{j+1})$$
  
= 
$$\sum_{j=1}^{d-1} (\Delta_j - \Delta_{j+1}) \sum_{i=1}^{j} w_i(m_i - y_i)$$
  
\ge 0.

In the latter step we used the assumptions on m and (5.17).

# 5.2.4 The pool-adjacent-violators algorithm (PAVA)

Now we consider a generalization of the function f in the previous subsection, namely,

$$f(\boldsymbol{m}) := \sum_{i=1}^d h_i(m_i)$$

with certain functions  $h_i : \mathbb{R} \mapsto (-\infty, \infty]$ . Explicit examples for the  $h_i$ , which are relevant in statistical applications, are:

- $h_i(r) := w_i(y_i r)^2$ ,
- $h_i(r) := |y_i r|,$
- $h_i(r) := w_i (y_i L(r) + (1 y_i) L(1 r))$ , where  $y_i \in [0, 1]$  and

$$L(r) := \begin{cases} -\log(s) & \text{if } s \in (0,1] \\ \infty & \text{else.} \end{cases}$$

Assumption on the functions  $h_i$ . For any nonvoid index set  $Q \subset \{1, 2, \dots, d\}$  let

$$h_Q(r) := \sum_{i \in Q} h_i(r).$$

We assume that there exists a number  $\xi_Q \in \mathbb{R}$  with the following properties:  $h_Q(\xi_Q) < \infty$ , and

(5.18) 
$$h_Q$$
 is   
 $\begin{cases} \text{isotonic on } [\xi_Q, \infty), \\ \text{antitonic on } (-\infty, \xi_Q] \end{cases}$ 

In the explicit examples mentioned before, this assumption is always satisfied:

- In case of  $h_i(r) = w_i(y_i r)^2$  or  $h_i(r) := w_i(y_iL(r) + (1 y_i)L(1 r)),$  $\xi_Q = \sum_{i \in Q} w_i y_i / \sum_{i \in Q} w_i.$
- In case of  $h_i(r) = |y_i r|$  we may choose

$$\xi_Q \in \text{Median}(y_i : i \in Q).$$

Abstract description of PAVA. In principle the algorithm operates on partitions  $\mathcal{P}$  of the index set  $\{1, 2, \ldots, d\}$  into blocks, that is, into sets of the form  $\{a, \ldots, b\}$  with integers  $1 \le a \le b \le d$ . For each such partition  $\mathcal{P}$  we define the vector

$$\boldsymbol{m}_{\mathcal{P}} := \left(\sum_{P \in \mathcal{P}} \mathbb{1}_{[i \in P]} \xi_P\right)_{i=1}^d \in \operatorname*{arg\,min}_{\boldsymbol{m} \in \mathcal{X}_{\mathcal{P}}} f(\boldsymbol{m}),$$

where

$$\mathcal{X}_{\mathcal{P}} := \left\{ \left( \sum_{P \in \mathcal{P}} 1_{[i \in P]} v_P \right)_{i=1}^d : v_P \in \mathbb{R} \text{ for } P \in \mathcal{P} \right\}$$

Initialization: We start with the finest partition

$$\mathcal{P} = \{\{1\}, \{2\}, \dots, \{n\}\}.$$

Induction step: Suppose that the vector  $m_{\mathcal{P}}$  is not in  $\mathcal{K}_{\uparrow}$ . More generally, suppose that the partition  $\mathcal{P}$  contains two neighboring blocks  $Q = \{a, \ldots, b-1\}$  and  $R = \{b, \ldots, c\}$  such that  $\xi_Q > \xi_R$ , that is, Q and R are two "adjacent violators". Then we replace  $\mathcal{P}$  with a new partition by replacing Q and R with their union  $Q \cup R = \{a, \ldots, c\}$  ("pool adjacent violators").

*End:* The induction step is repeated until finally  $m_{\mathcal{P}} \in \mathcal{K}_{\uparrow}$ . As shown next, this vector  $m_{\mathcal{P}}$  is automatically a minimizer of f over  $\mathcal{K}_{\uparrow}$ .

Validity of PAVA. Since the cardinality of  $\mathcal{P}$  decreases by one in each instance of the induction step, the algorithm will terminate after at most d-1 repetitions of the induction step. Obviously the final vector  $m_{\mathcal{P}}$  belongs to  $\mathcal{K}_{\uparrow}$ . That it minimizes f follows from the following observation:

**Lemma 5.25.** Let  $\mathcal{P}$  be a partition of  $\{1, 2, ..., d\}$  into blocks, among them the two neighbors  $Q = \{a, ..., b-1\}$  and  $R = \{b, ..., c\}$ . Further let  $\widetilde{\mathcal{P}}$  be the partition resulting from pooling Q and R. Suppose that  $\xi_Q \ge \xi_R$ . Then for each vector  $\mathbf{m} \in \mathcal{K}_{\uparrow} \cap \mathcal{X}_{\mathcal{P}}$  there exists a vector  $\widetilde{\mathbf{m}} \in \mathcal{K}_{\uparrow} \cap \mathcal{X}_{\widetilde{\mathcal{P}}}$  such that  $f(\widetilde{\mathbf{m}}) \le f(\mathbf{m})$ .

This lemma implies that the final partition  $\mathcal{P}$  produced by PAVA satisfies

$$\min_{\boldsymbol{m}\in\mathcal{K}_{\uparrow}} f(\boldsymbol{m}) = \min_{\boldsymbol{m}\in\mathcal{K}_{\uparrow}\cap\mathcal{X}_{\mathcal{P}}} f(\boldsymbol{m}) \geq \min_{\boldsymbol{m}\in\mathcal{X}_{\mathcal{P}}} f(\boldsymbol{m}) = f(\boldsymbol{m}_{\mathcal{P}}).$$

And since the vector  $m_{\mathcal{P}}$  belongs to  $\mathcal{K}_{\uparrow}$ , the previous inequality is even an equality, so  $m_{\mathcal{P}}$  minimizes f over  $\mathcal{K}_{\uparrow}$ .

**Proof of Lemma 5.25.** Let  $m \in \mathcal{K}_{\uparrow} \cap \mathcal{X}_{\mathcal{P}}$ , that means,

$$\boldsymbol{m} = \left(\sum_{P \in \mathcal{P}} 1_{[i \in P]} m_P\right)_{i=1}^d$$

with certain numbers  $m_P, P \in \mathcal{P}$  such that  $m_P \leq m_{P'}$  whenever  $\max(P) < \min(P')$ . Then

$$f(\boldsymbol{m}) = \sum_{P \in \mathcal{P}} h_P(m_P).$$

In case of  $\xi_Q > m_Q$  one could replace both numbers  $m_Q$  and  $m_R$  with the value

$$\min(\xi_Q, m_R) \in [m_Q, m_R]$$

without increasing  $f(\boldsymbol{m})$ , because  $h_Q$  is antitonic on  $(-\infty, \xi_Q]$  while  $h_R$  is isotonic on  $[\xi_R, \infty) \supset [\xi_Q, \infty)$ .

In case of  $\xi_Q \leq m_Q$  it follows from  $\xi_Q \geq \xi_R$  that  $\xi_R \leq m_Q$ , so we could replace  $m_R$  with  $m_Q$  without increasing  $f(\boldsymbol{m})$ .

In both cases the resulting vector belongs to the smaller set  $\mathcal{K}_{\uparrow} \cap \mathcal{X}_{\widetilde{\mathcal{P}}}$ .

**Explicit implementations.** In practice, it is even advisable to pool adjacent weak violators, that is, to pool two adjacent blocks  $Q = \{a, ..., b - 1\}$  and  $R = \{b, ..., c\}$  in  $\mathcal{P}$  whenever  $\xi_Q \ge \xi_R$ . This is justified by Lemma 5.25, too.

We introduce an auxiliary integer variable c running from 2 through d, and in each step we compute a vector in

$$\operatorname*{arg\,min}_{m{m}\in\mathbb{R}^c:m_1\leq\cdots\leq m_c}\sum_{i=1}^ch_i(m_i).$$

To this end we use an additional tuple  $\ell = (\ell_0, \ell_1, \dots, \ell_k)$  of variable length k + 1 with integer components  $0 = \ell_0 < \ell_1 < \ell_2 < \dots < \ell_k = c$ . This tuple corresponds to the partition

$$\mathcal{P} = \left(\underbrace{\{1,\ldots,\ell_1\}}_{=:P_1},\underbrace{\{\ell_1+1,\ldots,\ell_2\}}_{=:P_2},\ldots,\underbrace{\{\ell_{k-1}+1,\ldots,c\}}_{=:P_k}, \{c+1\},\{c+2\},\ldots,\{d\}\right)$$

in the abstract definition. In addition we use the tuple  $M = (M_0, M_1, \dots, M_k)$  with components  $M_0 := -\infty$  and  $M_j := \xi_{P_j}$  for  $1 \le j \le k$ . Then the algorithm works as in Table 5.3.

```
Algorithm \boldsymbol{m} \leftarrow \textbf{PAVA}(h_1, \dots, h_d)
\begin{array}{l} \boldsymbol{M} \leftarrow (-\infty, \xi_{\{1\}}) \\ \boldsymbol{k} \leftarrow 1 \end{array}
for c \leftarrow 2 to d do
             \boldsymbol{\ell} \leftarrow (\boldsymbol{\ell}, c)
             \boldsymbol{M} \leftarrow (\boldsymbol{M}, \xi_{\{c\}})
             k \leftarrow k+1
             while M_{k-1} \ge M_k do
                          \boldsymbol{\ell} \leftarrow (\ell_0, \dots, \ell_{k-2}, c)
                           M \leftarrow (M_0, \dots, M_{k-2}, \xi_{\{\ell_{k-2}+1, \dots, c\}})
                           k \leftarrow k - 1
             end while
end for
\boldsymbol{m} \leftarrow (0)_{i=1}^d
for j \leftarrow 1 to k do
             for i \leftarrow \ell_{j-1} + 1 to \ell_j do
                        m_i \leftarrow M_j
             end for
end for.
```

Table 5.3: Explicit version of PAVA.

• The special case  $h_i(r) = w_i(y_i - r)^2$ . Here one should use an additional tuple  $W = (W_1, \ldots, W_k)$  with the weights  $W_j := \sum_{i \in P_j} w_i$ . For then

$$\xi_{\{\ell_{k-2}+1,\dots,c\}} = \frac{W_{k-1}M_{k-1} + W_kM_k}{W_{k-1} + W_k} \quad \text{and} \quad \sum_{i \in P_{k-1} \cup P_k} w_i = W_{k-1} + W_k.$$

Thus each instance of the while–loop in Table 5.3 requires only O(1) operations. Consequently, the whole algorithm requires O(d) memory and O(d) operations.

• The special case  $h_i(r) = |y_i - r|$ . Here one should use an additional vector  $\mathbf{z} = (z_i)_{i=1}^d$ which is equal to  $\mathbf{y}$  initially. In general,  $(z_i)_{i \in P_j}$  contains the components of  $(y_i)_{i \in P_j}$  in nondecreasing order. If two blocks  $P_{k-1} = \{\ell_{k-1}, \ldots, \ell_k - 1\}$  and  $P_k = \{\ell_k, \ldots, c\}$  have to be pooled, the entries  $z_{\ell_{k-1}}, \ldots, z_c$  need to be sorted which is possible within O(d) steps (similarly as in "MergeSort"). All in all the whole algorithm requires O(d) memory and  $O(d^2)$  operations.

**Refinement.** Suppose that the local minimizers  $\xi_Q$ ,  $\emptyset \neq Q \subset \{1, \ldots, d\}$ , are unique. This implies that for arbitrary disjoint and nonempty sets  $Q, R \subset \{1, \ldots, d\}$ ,

$$\min(\xi_Q,\xi_R) \leq \xi_{Q\cup R} \leq \max(\xi_Q,\xi_R)$$

Consequently, instead of starting the PAVA (abstract description) with the finest possible partition  $\mathcal{P} = (\{1\}, \{2\}, \dots, \{d\})$ , one can start with a partition into maximal blocks on which  $i \mapsto \xi_{\{i\}}$  is non-increasing.

For the special case of  $h_i(r) = w_i(y_i - r)^2$ , Table 5.4 contains pseudocode for a refined version of the PAVA. Table 5.5 illustrates this procedure for a vector  $y \in \mathbb{R}^{12}$  and  $w = \mathbf{1}_{12}$ .

Algorithm $m{m} \leftarrow  extsf{PAVA}(m{y},m{w})$	
$\ell \leftarrow (0)$	% initialize
$\boldsymbol{W} \leftarrow (0)$	
$M \leftarrow (-\infty)$	
$c \leftarrow 0$	
$k \leftarrow 0$	
while $c < d$ do	
$c \leftarrow c + 1$	% add new block
$W_{\text{new}} \leftarrow w_c$	
$S_{ ext{new}} \leftarrow w_c y_c$	
$c' \leftarrow c+1$	
while $c' \leq d$ and $y_{c'} \leq y_c$ do	
$c \leftarrow c'$	
$W_{\text{new}} \leftarrow W_{\text{new}} + w_c$	
$S_{\text{new}} \leftarrow S_{\text{new}} + w_c y_c$	
$c' \leftarrow c+1$	
end while	
$\boldsymbol{\ell} \leftarrow (\boldsymbol{\ell},c)$	
$\boldsymbol{W} \leftarrow (\boldsymbol{W}, W_{\text{new}})$	
$\boldsymbol{M} \leftarrow (\boldsymbol{M}, S_{\text{new}} / W_{\text{new}})$	
$k \leftarrow k + 1$	
while $M_{k-1} \ge M_k$ do	
$\boldsymbol{\ell} \leftarrow (\ell_0, \dots, \ell_{k-2}, c)$	% pool adjacent
$W_{\text{new}} \leftarrow W_{k-1} + W_k$	% weak violators
$\boldsymbol{M} \leftarrow \left(M_0, \dots, M_{k-2}, (W_{k-1}M_{k-1} + W_kM_k)/W_{\text{new}}\right)$	
$\boldsymbol{W} \leftarrow \left(W_0, \dots, W_{k-1}, W_{\text{new}}\right)$	
$k \leftarrow k-1$	
end while	
end for	
$oldsymbol{m} \leftarrow (0)_{i=1}^d$	% generate m
for $j \in \{1, \dots, k\}$ do	
for $i \in \{\ell_{j-1}+1,\ldots,\ell_j\}$ do	
$m_i \leftarrow M_j$	
end for	
end for.	

Table 5.4: Explicit version of PAVA for isotonic weighted least squares.

The next three exercises show that the PAVA leads to an explicit nonparametric density estimator introduced by U. Grenander (1956). Instead of minimizing a convex function f over the cone  $\mathcal{K}_{\uparrow}$  one has to use an appropriate variant of the PAVA for the cone  $\mathbb{R}^n_{\downarrow} = \{ \boldsymbol{m} \in \mathbb{R}^n : m_1 \geq \cdots \geq m_n \}.$ 

**Exercise 5.26** (Silverman's trick). Let  $\mathcal{F}$  be a family of probability densities on a measure space  $(\mathcal{X}, \mathcal{B}, \mu)$ , and let  $\mathcal{P} := \{tf : t > 0, f \in \mathcal{F}\}$  be the cone generated by  $\mathcal{F}$ . Show that for given
	$\mid y^{ op}$	=		(	2	0	3	4	1	0	3	5	6	1	7	3	)
add	$\ell$	$\leftarrow$	(	0		2	)										
new	W	$\leftarrow$	(	0		2	)										
block	$\mid M$	$\leftarrow$	(	$-\infty$		1	)										
add	l	$\leftarrow$	(	0		2	3	)									
new	W	$\leftarrow$	(	0		2	1	)									
block	$\mid M$	$\leftarrow$	(	$-\infty$		1	3	)									
add	l	$\leftarrow$	(	0		2	3			6	)						
new	W	$\leftarrow$	(	0		2	1			3	)						
block	$\mid M$	$\leftarrow$	(	$-\infty$		1	3			$1.\overline{6}$	)						
pool	$\ell$	$\leftarrow$	(	0		2				6	)						
adj.	W	$\leftarrow$	(	0		2				4	)						
viol.	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	)						
add	l	$\leftarrow$	(	0		2				6	7	)					
new	W	$\leftarrow$	(	0		2				4	1	)					
block	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	3	)					
add	l	$\leftarrow$	(	0		2				6	7	8	)				
new	W	$\leftarrow$	(	0		2				4	1	1	)				
block	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	3	5	)				
add	$\ell$	$\leftarrow$	(	0		2				6	7	8		10	)		
new	W	$\leftarrow$	(	0		2				4	1	1		2	)		
block	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	3	5		3.5	)		
pool	l	$\leftarrow$	(	0		2				6	7			10	)		
adj.	W	$\leftarrow$	(	0		2				4	1			3	)		
viol.	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	3			4	)		
add	l	$\leftarrow$	(	0		2				6	7			10		12	)
new	W	$\leftarrow$	(	0		2				4	1			3		2	)
block	$\mid M$	$\leftarrow$	(	$-\infty$		1				2	3			4		5	)
	$m^ op$	=		(	1	1	2	2	2	2	3	4	4	4	5	5	)

Table 5.5: Numerical example for the PAVA with  $h_i(r) = (y_i - r)^2$ .

 $n \in \mathbb{N}$  and  $x_1, \ldots, x_n \in \mathcal{X}$ ,

$$\underset{f \in \mathcal{F}}{\operatorname{arg\,max}} \ \ell(f) \ = \ \underset{p \in \mathcal{P}}{\operatorname{arg\,min}} \Big( -\ell(p) + n \int p \, d\mu \Big)$$

where  $\ell(p) := \sum_{i=1}^{n} \log p(x_i)$ .

**Exercise 5.27** (Grenander's estimator, I). Let  $\mathcal{F}$  be the family of monotone decreasing probability densities f on  $(0, \infty)$  with respect to Lebesgue measure. Let  $0 = x_0 < x_1 < x_2 < \cdots x_n$  be given numbers. Viewing  $x_1, \ldots, x_n$  as the order statistics of n independent random variables with unknown density  $f \in \mathcal{F}$ , the nonparametric maximum likelihood estimator of f is defined as

$$\widehat{f} := rgmax_{f \in \mathcal{F}} \ell(f) \quad \text{with} \quad \ell(f) := \sum_{i=1}^n \log f(x_i).$$

The goal of this and the next exercise is to show that this estimator is well-defined and can be computed by means of the pool-adjacent-violators algorithm.

(a) Show that for each  $f \in \mathcal{F}$  there exists a density  $\overline{f} \in \mathcal{F}$  such that

$$\bar{f} \equiv \begin{cases} \bar{f}(x_i) & \operatorname{on}(x_{i-1}, x_i] \text{ for } 1 \le i \le n, \\ 0 & \operatorname{on}(x_n, \infty), \end{cases}$$

and  $\ell(\bar{f}) \ge \ell(f)$  with equality if and only if  $f \equiv \bar{f}$ .

(b) Show by means of the previous part and Exercise 5.26 that  $\hat{f}$  is uniquely determined by

$$\widehat{\boldsymbol{p}} := \operatorname*{arg\,min}_{\boldsymbol{p} \in \mathbb{R}^n_{\downarrow}} \sum_{i=1}^n h_i(p_i) \quad \text{with} \quad h_i(s) := -\log(s) + n(x_i - x_{i-1})s$$

**Exercise 5.28** (Grenander's estimator, II). (a) Show that the vector  $\hat{p}$  in Exercise 5.27 is identical with

$$\underset{m \in \mathbb{R}^n_{\downarrow}}{\operatorname{arg\,min}} \sum_{i=1}^n w_i (y_i - m_i)^2 \quad \text{with} \quad w_i := x_i - x_{i-1} \text{ and } y_i := (nw_i)^{-1}.$$

(b) Let  $\widehat{F}_{emp} : [0, \infty) \to [0, 1]$  be the empirical distribution function of the points  $x_1, \ldots, x_n$ , that means,  $\widehat{F}_{emp}(x) = i/n$  for  $x_i \leq x < x_{i+1}$ ,  $0 \leq i \leq n$ , with  $x_{n+1} := \infty$ . Further let  $\widehat{F}(x) := \int_0^x \widehat{f}(t) dt$  with Grenander's estimator  $\widehat{f}$ . Verify that  $\widehat{F}$  is the smallest concave function which is bounded from below by  $\widehat{F}_{emp}$ .



Figure 5.2: Example for concave regression. Upper panel: data pairs  $(X_{\ell}, Y_{\ell})$  ( $\circ$ ) and preprocessed pairs  $(t_i, y_i)$  ( $\bullet$ ). Lower panel: data pairs and final fit.



Figure 5.3: Example for concave regression: Some intermediate steps of the active set algorithm.



Figure 5.4: Isotonic regression and convex minorants.

## **Chapter 6**

# **Conjugate Gradients**

The main part of this chapter is based on the textbook of G. Opfer (1994).

#### 6.1 The Task

Let  $(V, \langle \cdot, \cdot \rangle, \|\cdot\|)$  be a Euclidean space<sup>1</sup> with dimension d, and let  $A : V \to V$  be a self-adjoint, positive definite linear operator. That A is self-adjoint means that

$$\langle \boldsymbol{x}, \boldsymbol{A} \boldsymbol{y} \rangle = \langle \boldsymbol{A} \boldsymbol{x}, \boldsymbol{y} \rangle$$
 for arbitrary  $\boldsymbol{x}, \boldsymbol{y} \in \boldsymbol{V}$ .

And a self-adjoint linear operator is called positive definite if

$$\langle \boldsymbol{x}, \boldsymbol{A} \boldsymbol{x} \rangle > 0$$
 for arbitrary  $\boldsymbol{x} \in \boldsymbol{V} \setminus \{\boldsymbol{0}\}$ .

In particular, A is invertible. Now, for a given vector  $b \in V$ , our goal is to determine the vector

$$\boldsymbol{x}_* := \boldsymbol{A}^{-1} \boldsymbol{b}.$$

This looks like a rather trivial task. If we choose an orthonormal basis  $u_1, u_2, \ldots, u_d$  of V, then A corresponds to a symmetric, positive definite matrix in  $\mathbb{R}^{d \times d}$ , b may be viewed as a vector in  $\mathbb{R}^d$ , and  $x_*$  is the solution of the linear equation system Ax = b. The reason for considering an operator rather than a matrix is that in some applications it is easy to compute the point Ax for any  $x \in V$ , but the storage or inversion of the corresponding matrix A would require too much space or computation time. Then it is not obvious how to determine  $x_*$ , and in what follows, a special iterative scheme will be derived.

Note first that the vector  $x_*$  is the unique minimizer of the function  $f: V \to \mathbb{R}$  with

$$f(\boldsymbol{x}) \ := \ 2^{-1} \langle \boldsymbol{x}, \boldsymbol{A} \boldsymbol{x} 
angle - \langle \boldsymbol{b}, \boldsymbol{x} 
angle_{+}$$

Indeed, for arbitrary  $x, h \in V$ ,

(6.1) 
$$f(\boldsymbol{x}+\boldsymbol{h}) = f(\boldsymbol{x}) - \langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{h} \rangle + 2^{-1} \langle \boldsymbol{h}, \boldsymbol{A} \boldsymbol{h} \rangle$$

<sup>&</sup>lt;sup>1</sup>a finite-dimensional real Hilbert space

with the "residual"

$$r(x) := b - Ax = A(x_* - x)$$

This expansion shows that  $\nabla f(x) = -r(x)$  and  $D^2 f(x) = A$ . In particular,

$$f(\boldsymbol{x}_* + \boldsymbol{h}) = f(\boldsymbol{x}_*) + 2^{-1} \langle \boldsymbol{h}, \boldsymbol{A} \boldsymbol{h} \rangle \geq f(\boldsymbol{x}_*)$$

with equality if and only if h = 0. Precisely, for arbitrary  $x \in V$ ,

$$f(\boldsymbol{x}) - f(\boldsymbol{x}_*) = 2^{-1} \langle \boldsymbol{x} - \boldsymbol{x}_*, \boldsymbol{A}(\boldsymbol{x} - \boldsymbol{x}_*) \rangle$$

(6.2) 
$$= 2^{-1} \langle \boldsymbol{A}^{-1} \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{r}(\boldsymbol{x}) \rangle$$

(6.3) 
$$= 2^{-1} \| \boldsymbol{A}^{1/2} (\boldsymbol{x} - \boldsymbol{x}_*) \|^2.$$

#### 6.2 The Gradient Method

A first approach is to minimize f iteratively in the direction of its residual. That means, for a vector  $x \in V \setminus \{x_*\}$  we define

$$\boldsymbol{\psi}(\boldsymbol{x}) := \boldsymbol{x} + t(\boldsymbol{x})\boldsymbol{r}(\boldsymbol{x})$$

with

$$t(oldsymbol{x}) \ := \ rgmin_{t \in \mathbb{R}} f(oldsymbol{x} + toldsymbol{r}(oldsymbol{x})).$$

Starting from an arbitrary point  $x_0 \in V$ , we define the sequence  $(x_k)_{k\geq 0}$  inductively via  $x_{k+1} := \psi(x_k)$ , where  $\psi(x_*) := x_*$ . As shown later, this sequence will always converge to the minimizer  $x_*$ .

Since

$$f(\boldsymbol{x}+t\boldsymbol{r}(\boldsymbol{x})) = f(\boldsymbol{x}) - t \|\boldsymbol{r}(\boldsymbol{x})\|^2 + 2^{-1}t^2 \langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}\boldsymbol{r}(\boldsymbol{x}) \rangle,$$

the number t(x) is given by

$$t(oldsymbol{x}) \;=\; rac{\|oldsymbol{r}(oldsymbol{x})\|^2}{\langleoldsymbol{r}(oldsymbol{x}),oldsymbol{A}oldsymbol{r}(oldsymbol{x})
angle},$$

and

$$f(\boldsymbol{\psi}(\boldsymbol{x})) = f(\boldsymbol{x}) - \frac{\|\boldsymbol{r}(\boldsymbol{x})\|^4}{2\langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}\boldsymbol{r}(\boldsymbol{x}) \rangle}$$

Combining this equation and (6.2) leads to

$$\frac{f(\boldsymbol{\psi}(\boldsymbol{x})) - f(\boldsymbol{x}_*)}{f(\boldsymbol{x}) - f(\boldsymbol{x}_*)} \; = \; 1 - \frac{\|\boldsymbol{r}(\boldsymbol{x})\|^4}{\langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}\boldsymbol{r}(\boldsymbol{x})\rangle \langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}^{-1}\boldsymbol{r}(\boldsymbol{x})\rangle}$$

The ratio on the right hand side may be bounded as follows.

**Lemma 6.1.** Let  $Ax = \sum_{i=1}^{d} \lambda_i \langle u_i, x \rangle u_i$  for  $x \in V$  with an orthonormal basis  $u_1, u_2, \ldots, u_d$  of V and eigenvalues  $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_d > 0$ . For arbitrary  $y \in V \setminus \{0\}$ ,

$$1 \leq rac{\langle oldsymbol{y},oldsymbol{A}oldsymbol{y},oldsymbol{A}^{-1}oldsymbol{y}
angle}{\|oldsymbol{y}\|^4} \leq rac{(\lambda_1+\lambda_d)^2}{4\lambda_1\lambda_d}.$$

The second last inequality is an equality if and only if y is an eigenvector of A. The last inequality is an equality if y is a multiple of  $u_1 + u_d$  or  $u_1 - u_d$ .

For the algorithmic mapping  $\psi$ , this lemma leads to the following inequality:

**Theorem 6.2**. With the same notation as in Lemma 6.1,

$$rac{f(oldsymbol{\psi}(oldsymbol{x})) - f(oldsymbol{x}_*)}{f(oldsymbol{x}) - f(oldsymbol{x}_*)} \ \le \ \Big(rac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d}\Big)^2$$

for arbitrary  $x \in V \setminus \{x_*\}$ . Equality holds if r(x) is a multiple of  $u_1 + u_d$  or  $u_1 - u_d$ . In that case,  $r(\psi(x))$  is a multiple of  $u_1 - u_d$  or  $u_1 + u_d$ , respectively.

This theorem implies that the sequence  $(x_k)_{k\geq 0}$  just introduced converges to  $x_*$  for any starting point  $x_0 \in V$ . More precisely, the inequality is equivalent to

$$\frac{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{\psi}(\boldsymbol{x})-\boldsymbol{x}_*)\right\|}{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{x}-\boldsymbol{x}_*)\right\|} \ \le \ \frac{\lambda_1-\lambda_d}{\lambda_1+\lambda_d} \ = \ \frac{1-\lambda_d/\lambda_1}{1+\lambda_d/\lambda_1} \ = \ 1-\frac{2}{\lambda_1/\lambda_d+1}$$

by virtue of (6.3). Hence,

$$ig\Vert oldsymbol{A}^{1/2}(oldsymbol{x}_k-oldsymbol{x}_*)ig\Vert\ \leq\ \Big(rac{1-\lambda_d/\lambda_1}{1+\lambda_d/\lambda_1}\Big)^k ig\Vert oldsymbol{A}^{1/2}(oldsymbol{x}_0-oldsymbol{x}_*)ig\Vert$$

for arbitrary  $k \ge 0$  with equality in case of  $r(x_0)$  being a multiple of  $u_1 \pm u_d$ . However, the speed of convergence can be rather low whenever the condition number  $\lambda_d/\lambda_1$  of the operator A is close to zero.

Figure 6.1 illustrates the gradient method in case of  $V = \mathbb{R}^2$  with the usual inner product, A = diag(1,3),  $b = (1,1)^{\top}$  whence  $x_* = (5,5/3)^{\top}$ , and  $x_0 = 0$ . Here  $r(x_k)$  is always a positive multiple of  $(1,1)^{\top}$  or  $(1,-1)^{\top}$ , and  $f(x_k) - f(x_*)$  decreases in each step by the same factor  $((3-1)/(3+1))^2 = 1/4$ .



Figure 6.1: The gradient method.

**Proof of Lemma 6.1.** The asserted inequalities are equivalent to the following ones: For arbitrary vectors  $p \in [0, 1]^d$  with  $\sum_{i=1}^d p_i = 1$ ,

$$1 \leq \sum_{i=1}^{d} p_i \lambda_i \sum_{j=1}^{d} p_j \lambda_j^{-1} \leq \frac{(\lambda_1 + \lambda_d)^2}{4\lambda_1 \lambda_d}.$$

Indeed, if we write  $oldsymbol{y} = \sum_{i=1}^d a_i oldsymbol{u}_i$  with  $oldsymbol{a} \in \mathbb{R}^d$ , then

$$\|\boldsymbol{y}\|^{-2}\langle \boldsymbol{y}, \boldsymbol{A}^{s}\boldsymbol{y}\rangle = \sum_{i=1}^{d} p_{i}\lambda_{i}^{s} \text{ with } p_{i} := \|\boldsymbol{y}\|^{-2}a_{i}^{2}.$$

Furthermore, any probability vector p may be obtained in this fashion from some vector  $a \in \mathbb{R}^d \setminus \{0\}$ .

In the center of the new inequalities we see  $\mathbb{E}(X) \mathbb{E}(X^{-1})$  with a random variable X > 0 having distribution  $\sum_{i=1}^{d} p_i \delta_{\lambda_i}$ . Since the function  $x \mapsto x^{-1}$  is strictly convex on  $(0, \infty)$ , Jensen's inequality implies that  $\mathbb{E}(X^{-1}) \ge \mathbb{E}(X)^{-1}$ . Equality holds if and only if X has a degenerate distribution. In the original inequality, this means that y is an eigenvector of A.

On the other hand,  $\lambda_d \leq X \leq \lambda_1$ , so convexity of  $x \mapsto x^{-1}$  yields the following inequality:

$$X^{-1} \leq \frac{\lambda_1 - X}{\lambda_1 - \lambda_d} \lambda_d^{-1} + \frac{X - \lambda_d}{\lambda_1 - \lambda_d} \lambda_1^{-1}.$$

Consequently,

$$\mathbb{E}(X) \mathbb{E}(X^{-1}) \leq \mathbb{E}(X) \left( \frac{\lambda_1 - \mathbb{E}(X)}{\lambda_1 - \lambda_d} \lambda_d^{-1} + \frac{\mathbb{E}(X) - \lambda_d}{\lambda_1 - \lambda_d} \lambda_1^{-1} \right) \\
= \frac{\mathbb{E}(X)}{(\lambda_1 - \lambda_d)\lambda_1\lambda_d} \left( \lambda_1^2 - \lambda_d^2 - (\lambda_1 - \lambda_d) \mathbb{E}(X) \right) \\
= \frac{\mathbb{E}(X)(\lambda_1 + \lambda_d - \mathbb{E}(X))}{\lambda_1\lambda_d} \\
\leq \frac{(\lambda_1 + \lambda_d)^2}{4\lambda_1\lambda_d}.$$

Equality holds if  $p_1 = p_d = 1/2$ , and in the original inequality, this means that  $\boldsymbol{y}$  is proportional to  $\boldsymbol{u}_1 \pm \boldsymbol{u}_d$ .

#### Proof of Theorem 6.2. According to Lemma 6.1,

$$\frac{f(\boldsymbol{\psi}(\boldsymbol{x})) - f(\boldsymbol{x}_*)}{f(\boldsymbol{x}) - f(\boldsymbol{x}_*)} = 1 - \frac{\|\boldsymbol{r}(\boldsymbol{x})\|^4}{\langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}\boldsymbol{r}(\boldsymbol{x}) \rangle \langle \boldsymbol{r}(\boldsymbol{x}), \boldsymbol{A}^{-1}\boldsymbol{r}(\boldsymbol{x}) \rangle} \leq 1 - \frac{4\lambda_1\lambda_d}{(\lambda_1 + \lambda_d)^2} = \left(\frac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d}\right)^2.$$

In general,

$$\boldsymbol{x}_* - \boldsymbol{\psi}(\boldsymbol{x}) = \boldsymbol{x}_* - \boldsymbol{x} - t(\boldsymbol{x})\boldsymbol{r}(\boldsymbol{x}) = (\boldsymbol{A}^{-1} - t(\boldsymbol{x})\boldsymbol{I})\boldsymbol{r}(\boldsymbol{x}).$$

Consequently, if  $\boldsymbol{r}(\boldsymbol{x}) = \gamma(\boldsymbol{u}_1 \pm \boldsymbol{u}_d)$  for some  $\gamma \neq 0$ , then

$$egin{aligned} m{r}(\psi(m{x})) &= m{A}(m{x}_* - \psi(m{x})) &= (m{I} - t(m{x})m{A})m{r}(m{x}) \ &= igg(m{I} - rac{\|m{r}(m{x})\|^2}{\langlem{r}(m{x}),m{A}m{r}(m{x})
angle}m{A}igg)m{r}(m{x}) \ &= igg(m{u}_1 \pm m{u}_d - rac{2}{\lambda_1 + \lambda_d}(\lambda_1m{u}_1 \pm \lambda_dm{u}_d)igg) \ &= igg(egin{aligned} &\lambda_d - \lambda_1 \\ \lambda_1 + \lambda_d \end{matrix} m{u}_1 \pm rac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d} m{u}_d igg) \ &= -rac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d} \gamma(m{u}_1 \mp m{u}_d). \end{aligned}$$

154

### 6.3 Conjugate Directions

As we shall see soon, the minimization of f would be rather simple if we had a basis of V which is "conjugate with respect to A".

**Definition 6.3** (Conjugate vectors). Vectors  $h_1, h_2, \ldots, h_k \in V$  are called *conjugate with respect to* A if

$$\langle \boldsymbol{h}_i, \boldsymbol{A}\boldsymbol{h}_j \rangle = 0 \quad \text{for } 1 \leq i < j \leq k$$

In other words, the vectors  $A^{1/2}h_1, A^{1/2}h_2, \dots, A^{1/2}h_k$  are pairwise orthogonal.

**Lemma 6.4.** Let  $h_1, h_2, \ldots, h_\ell \in V \setminus \{0\}$  be conjugate with respect to A. For an arbitrary point  $x_0 \in \mathbb{R}^d$  and any integer  $k \in \{1, 2, \ldots, \ell\}$  let

$$egin{array}{lll} m{V}_k &:= ext{span}(m{h}_1,\ldots,m{h}_k), \ m{x}_k &:= ext{argmin}_{m{x}\,\in\,m{x}_0+m{V}_k}\,f(m{x}). \end{array}$$

Then

$$oldsymbol{x}_k \;=\; oldsymbol{x}_0 + \sum_{i=1}^k t_i oldsymbol{h}_i \;=\; oldsymbol{x}_{k-1} + t_k oldsymbol{h}_k$$

with

$$t_i \ := \ rac{\langle m{r}(m{x}_0),m{h}_i
angle}{\langlem{h}_i,m{A}m{h}_i
angle} \ = \ rac{\langlem{r}(m{x}_{i-1}),m{h}_i
angle}{\langlem{h}_i,m{A}m{h}_i
angle}.$$

Moreover,  $\boldsymbol{r}(\boldsymbol{x}_k) \perp \boldsymbol{V}_k$ .

**Proof of Lemma 6.4.** For  $t \in \mathbb{R}^k$ , it follows from (6.1) that

$$f\left(\boldsymbol{x}_{0} + \sum_{i=1}^{k} t_{i}\boldsymbol{h}_{i}\right) = f(\boldsymbol{x}_{0}) - \sum_{i=1}^{k} t_{i}\langle \boldsymbol{r}(\boldsymbol{x}_{0}), \boldsymbol{h}_{i}\rangle + 2^{-1}\sum_{i,j=1}^{k} t_{i}t_{j}\langle \boldsymbol{h}_{i}, \boldsymbol{A}\boldsymbol{h}_{j}\rangle$$
$$= f(\boldsymbol{x}_{0}) + \sum_{i=1}^{k} \left(2^{-1}t_{i}^{2}\langle \boldsymbol{h}_{i}, \boldsymbol{A}\boldsymbol{h}_{i}\rangle - t_{i}\langle \boldsymbol{r}(\boldsymbol{x}_{0}), \boldsymbol{h}_{i}\rangle\right).$$

Obviously, this is minimal in t if and only if

$$t_i = rac{\langle \boldsymbol{r}(\boldsymbol{x}_0), \boldsymbol{h}_i 
angle}{\langle \boldsymbol{h}_i, \boldsymbol{A} \boldsymbol{h}_i 
angle}.$$

In case of  $i \geq 2$ ,

$$r(x_{i-1}) = b - Ax_{i-1} = b - Ax_0 + A(x_0 - x_{i-1}) = r(x_0) + A(x_0 - x_{i-1}),$$

and this implies that

$$\langle \boldsymbol{r}(\boldsymbol{x}_0), \boldsymbol{h}_i 
angle \; = \; \langle \boldsymbol{r}(\boldsymbol{x}_{i-1}), \boldsymbol{h}_i 
angle + \langle \boldsymbol{x}_{i-1} - \boldsymbol{x}_0, \boldsymbol{A} \boldsymbol{h}_i 
angle \; = \; \langle \boldsymbol{r}(\boldsymbol{x}_{i-1}), \boldsymbol{h}_i 
angle$$

because  $x_{i-1} - x_0 \in \text{span}(h_1, \dots, h_{i-1})$  and  $\langle h_j, Ah_i \rangle = 0$  for  $1 \leq j < i$ . This proves the alternative expression for the optimal  $t_i$ .

That  $r(x_k)$  is perpendicular to  $V_k$  follows from the fact that for arbitrary  $h \in V_k$ ,

$$0 = \frac{d}{dt}\Big|_{t=0} f(\boldsymbol{x}_k + t\boldsymbol{h}) = -\langle \boldsymbol{r}(\boldsymbol{x}_k), \boldsymbol{h} \rangle.$$

A Gram-Schmidt type procedure. The preceding considerations suggest the following procedure: We start with a vector  $x_0 \in V$  such that  $r(x_0) \neq 0$  which is equivalent to  $x_0 \neq x_*$ . Then we set

$$h_1 := r(x_0).$$

Suppose that for some  $k \in \mathbb{N}$  we have chosen conjugate vectors  $h_1, \ldots, h_k \in V \setminus \{0\}$  with respect to A. Let

$$oldsymbol{V}_k := \operatorname{span}(oldsymbol{h}_1,\ldots,oldsymbol{h}_k) \quad ext{and} \quad oldsymbol{x}_k := rgmin_{oldsymbol{x}} \inf_{oldsymbol{x} \in oldsymbol{x}_0 + oldsymbol{V}_k} f(oldsymbol{x}).$$

If  $r(x_k) = 0$ , we know that  $x_k = x_*$ . Otherwise we define

$$oldsymbol{h}_{k+1} \ := \ oldsymbol{r}(oldsymbol{x}_k) - \sum_{i=1}^k eta_{k,i}oldsymbol{h}_i$$

with certain coefficients  $\beta_{k,i} \in \mathbb{R}$  yet to be specified. Since  $r(x_k) \perp V_k$ , the vector  $h_{k+1}$  is nonzero for any choice of the  $\beta_{k,i}$ . The vectors  $h_1, \ldots, h_{k+1}$  are conjugate with respect to A if and only if for  $1 \leq j \leq k$ ,

$$0 = \langle \boldsymbol{h}_{k+1}, \boldsymbol{A}\boldsymbol{h}_j \rangle = \langle \boldsymbol{r}(\boldsymbol{x}_k), \boldsymbol{A}\boldsymbol{h}_j \rangle - \beta_{k,j} \langle \boldsymbol{h}_j, \boldsymbol{A}\boldsymbol{h}_j \rangle_{j}$$

that means,

$$eta_{k,j} \;=\; rac{\langle m{r}(m{x}_k),m{A}m{h}_j
angle}{\langlem{h}_j,m{A}m{h}_j
angle}$$

All in all, if we start from a point  $x_0 \neq x_*$ , we obtain linearly independent vectors  $h_1, \ldots, h_{k_o} \in V$  and points  $x_1, \ldots, x_{k_o} \in V$  such that for  $1 \leq k \leq k_o$ ,

$$oldsymbol{V}_k := \operatorname{span}(oldsymbol{h}_1,\ldots,oldsymbol{h}_k) = \operatorname{span}ig(oldsymbol{r}(oldsymbol{x}_0),\ldots,oldsymbol{r}(oldsymbol{x}_{k-1})ig) \quad ext{and} \quad oldsymbol{x}_k = rgmin_{oldsymbol{x}_0} \min_{oldsymbol{x}\inoldsymbol{x}_0+oldsymbol{V}_k} f(oldsymbol{x}).$$

Moreover, for  $1 \le k < k_o$ ,

$$oldsymbol{r}(oldsymbol{x}_k) \perp oldsymbol{V}_k \hspace{0.1in} ext{ and } \hspace{0.1in} oldsymbol{h}_{k+1} egin{cases} \in oldsymbol{r}(oldsymbol{x}_k) + oldsymbol{V}_k, \ oldsymbol{\perp} oldsymbol{A}oldsymbol{V}_k, \ \end{pmatrix}$$

while  $\boldsymbol{r}(\boldsymbol{x}_{k_o}) = \boldsymbol{0}$  and

$$\boldsymbol{x}_{k_o} = \boldsymbol{x}_{*}.$$

#### 6.4 The Conjugate Gradient (CG) Algorithm

At first glance, it looks as if we have solved the original problem already. But recall that we have situations in mind in which the storage of a matrix for A would be too involved. But the Gram-Schmidt type procedure just described would require to keep track of all vectors  $h_1, h_2, \ldots, h_{k_o}$ . Fortunately, as shown in the next theorem, all coefficients  $\beta_{k,j}$  with  $1 \le j < k$  are zero, so we only need  $h_k$  and  $r(x_k)$  to compute the next direction  $h_{k+1}$ .

**Theorem 6.5.** Starting from  $x_0 \neq x_*$ , let  $h_1, \ldots, h_{k_o}$  and  $x_1, \ldots, x_{k_o}$  be constructed as just described. For  $1 \leq k \leq k_o$ , the space  $V_k = \operatorname{span}(h_1, \ldots, h_k) = \operatorname{span}(r(x_0), \ldots, r(x_{k-1}))$  is equal to the so-called Krylow space

$$\operatorname{span}(\boldsymbol{A}^{i-1}\boldsymbol{r}(\boldsymbol{x}_0): 1 \le i \le k).$$

Moreover, for  $1 \le k \le k_o$ ,

$$m{x}_k = m{x}_{k-1} + t_k m{h}_k$$
 and  $m{r}(m{x}_k) = m{r}(m{x}_{k-1}) - t_k m{A} m{h}_k$  with  $t_k := rac{\|m{r}(m{x}_{k-1})\|^2}{\langlem{h}_k,m{A} m{h}_k 
angle}$ 

while in case of  $k < k_o$ ,

$$m{h}_{k+1} \;=\; m{r}(m{x}_k) + rac{\|m{r}(m{x}_k)\|^2}{\|m{r}(m{x}_{k-1})\|^2}\,m{h}_k$$

This theorem shows that  $h_{k+1} = r(x_k) - \sum_{i=1}^k \beta_{k,i} h_i$  is just a linear combination of  $r(x_k)$  and  $h_k$ . The explicit formulae in Theorem 6.5 yield a relatively simple algorithm for the computation (approximation) of  $x_*$ :

$\fbox{Algorithmus} \hspace{0.1cm} \boldsymbol{x}_k \leftarrow \hspace{0.1cm} \mathbf{CG}(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{x}_0, \epsilon)$
$ig  oldsymbol{r}_0 \leftarrow oldsymbol{b} - oldsymbol{A} oldsymbol{x}_0$
$oldsymbol{h}_1 \leftarrow oldsymbol{r}_0$
$k \leftarrow 0$
$\mathbf{while} \; \  \boldsymbol{r}_k \  > \epsilon \; \mathbf{do}$
$k \leftarrow k + 1$
$t_k \leftarrow \ oldsymbol{r}_{k-1}\ ^2 / \langleoldsymbol{h}_k,oldsymbol{A}oldsymbol{h}_k angle$
$oldsymbol{x}_k \leftarrow oldsymbol{x}_{k-1} + t_k oldsymbol{h}_k$
$oldsymbol{r}_k \leftarrow oldsymbol{r}_{k-1} - t_k oldsymbol{A} oldsymbol{h}_k$
$egin{array}{lll} eta_k \leftarrow \ m{r}_k\ ^2 / \ m{r}_{k-1}\ ^2 \end{array}$
$egin{array}{lll} m{h}_{k+1} \leftarrow m{r}_k + eta_km{h}_k \end{array}$
end while

After a finite number of steps, this algorithm yields a vector  $x_k$  such that  $||Ax_k - b|| \le \epsilon$ , where  $\epsilon \ge 0$  is an arbitrary given number. Here is a different description:

$$\begin{array}{l} \textbf{Algorithmus} \ \boldsymbol{x} \leftarrow \textbf{CG}(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{x}_{0}, \epsilon) \\ \boldsymbol{x} \leftarrow \boldsymbol{x}_{0} \\ \boldsymbol{r} \leftarrow \boldsymbol{b} - \boldsymbol{A} \boldsymbol{x} \\ \gamma \leftarrow \|\boldsymbol{r}\|^{2} \\ \boldsymbol{h} \leftarrow \boldsymbol{r} \\ \textbf{while} \ \sqrt{\gamma} > \epsilon \ \textbf{do} \\ \boldsymbol{y} \leftarrow \boldsymbol{A} \boldsymbol{h} \\ t \leftarrow \gamma/\langle \boldsymbol{h}, \boldsymbol{y} \rangle \\ \boldsymbol{x} \leftarrow \boldsymbol{x} + t \boldsymbol{h} \\ \boldsymbol{r} \leftarrow \boldsymbol{r} - t \boldsymbol{y} \\ \gamma_{\text{new}} \leftarrow \|\boldsymbol{r}\|^{2} \\ \boldsymbol{h} \leftarrow \boldsymbol{r} + (\gamma_{\text{new}}/\gamma) \boldsymbol{h} \\ \gamma \leftarrow \gamma_{\text{new}} \end{array}$$

$$\begin{array}{l} \textbf{end while} \end{array}$$

**Proof of Theorem 6.5.** We write  $r_{\ell} := r(x_{\ell})$  and  $W_k := \operatorname{span}(A^{i-1}r_0 : 1 \le i \le k)$ . Since  $\dim(V_k) = k$  and  $\dim(W_k) \le k$ , it suffices to show that  $V_k \subset W_k$  for  $1 \le k \le k_o$ . Obviously,  $W_1 = V_1$ . Suppose that  $V_k \subset W_k$  for some  $1 \le k < k_o$ . Then there exist real coefficients  $\lambda_1, \ldots, \lambda_k$  such that

$$oldsymbol{x}_k \ = \ oldsymbol{x}_0 + \sum_{i=1}^k \lambda_i oldsymbol{A}^{i-1} oldsymbol{r}_0$$

But then

$$m{r}_k \ = \ m{b} - m{A} m{x}_k \ = \ m{r}_0 - \sum_{i=1}^k \lambda_i m{A}^i m{r}_0 \ \in \ m{W}_{k+1}$$

whence  $V_{k+1} \subset W_{k+1}$ .

It remains to verify the recursion formulae for  $1 \le k \le k_o$ . We know already from Lemma 6.4 that

$$oldsymbol{x}_k \ = \ oldsymbol{x}_{k-1} + t_k oldsymbol{h}_k \quad ext{with} \quad t_k \ = \ rac{\langle oldsymbol{r}_{k-1}, oldsymbol{h}_k 
angle}{\langle oldsymbol{h}_k, oldsymbol{A} oldsymbol{h}_k 
angle}$$

But  $\boldsymbol{h}_k \in \boldsymbol{r}_{k-1} + \boldsymbol{V}_{k-1}$  and  $\boldsymbol{r}_{k-1} \perp \boldsymbol{V}_{k-1}$ , whence  $\langle \boldsymbol{r}_{k-1}, \boldsymbol{h}_k \rangle = \langle \boldsymbol{r}_{k-1}, \boldsymbol{r}_{k-1} \rangle = \| \boldsymbol{r}_{k-1} \|^2$ , and

$$t_k = \frac{\|\boldsymbol{r}_{k-1}\|^2}{\langle \boldsymbol{h}_k, \boldsymbol{A} \boldsymbol{h}_k \rangle}$$

Now the simplified formula for  $r_k$  follows from

$$r_k = b - Ax_k = b - Ax_{k-1} - A(x_k - x_{k-1}) = r_{k-1} - t_k Ah_k.$$

Concerning  $m{h}_{k+1}$ , recall that  $m{h}_{k+1} = m{r}_k - \sum_{i=1}^k eta_{k,i} m{h}_i$  with

$$eta_{k,i} = rac{\langle m{r}_k, m{A}m{h}_i 
angle}{\langle m{h}_i, m{A}m{h}_i 
angle}.$$

But for  $1 \leq i < k$ , the vector  $h_i$  lies in  $V_i = \operatorname{span}(A^{s-1}r_0 : 1 \leq s \leq i)$ , whence  $Ah_i \in \operatorname{span}(A^sr_0 : 1 \leq s \leq i) \subset V_k \perp r_k$ , and thus  $\beta_{k,i} = 0$ . Consequently, we know that

$$oldsymbol{h}_{k+1} \;=\; oldsymbol{r}_k - rac{\langle oldsymbol{r}_k, oldsymbol{A}oldsymbol{h}_k
angle}{\langle oldsymbol{h}_k, oldsymbol{A}oldsymbol{h}_k
angle} oldsymbol{h}_k.$$

Finally, it follows from  $\boldsymbol{r}_k = \boldsymbol{r}_{k-1} - t_k \boldsymbol{A} \boldsymbol{h}_k$  that

$$\boldsymbol{A}\boldsymbol{h}_{k} = t_{k}^{-1}(r_{k-1}-r_{k}) = \frac{\langle \boldsymbol{h}_{k}, \boldsymbol{A}\boldsymbol{h}_{k} \rangle}{\|\boldsymbol{r}_{k-1}\|^{2}}(\boldsymbol{r}_{k-1}-\boldsymbol{r}_{k}),$$

and this leads to

$$m{h}_{k+1} \;=\; m{r}_k - rac{\langle m{r}_k, m{r}_{k-1} - m{r}_k 
angle}{\|m{r}_{k-1}\|^2} \,m{h}_k \;=\; m{r}_k + rac{\|m{r}_k\|^2}{\|m{r}_{k-1}\|^2} \,m{h}_k,$$

because  $r_{k-1} \in V_k \perp r_k$ .

#### 6.5 Bounding the running time and approximation error

In many applications of the CG algorithm, it turns out that the norm of the residual  $r_k = b - Ax_k$ is rather small or even zero after  $k \ll d$  steps. This is related to the aforementioned Krylow spaces

$$\boldsymbol{V}_k = \operatorname{span}(\boldsymbol{A}^{i-1}\boldsymbol{r}_0: 1 \le i \le k)$$

with  $\boldsymbol{r}_0 = \boldsymbol{r}(\boldsymbol{x}_0)$ . As shown already,

$$\begin{split} \|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_k)\|^2 &= \min_{\boldsymbol{x} \in \boldsymbol{x}_0 + \boldsymbol{V}_k} \|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x})\|^2 \\ &= \min_{\boldsymbol{v} \in \boldsymbol{V}_k} \|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_0) - \boldsymbol{A}^{1/2}\boldsymbol{v}\|^2 \\ &= \min_{\boldsymbol{\beta} \in \mathbb{R}^k} \|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_0) - \boldsymbol{A}^{1/2}\sum_{j=1}^k \beta_j \boldsymbol{A}^{j-1}\boldsymbol{r}_0\|^2 \\ &= \min_{\boldsymbol{\beta} \in \mathbb{R}^k} \|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_0) - \sum_{j=1}^k \beta_j \boldsymbol{A}^j \boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_0)\|^2. \end{split}$$

Now we consider a spectral representation of A, that means,  $Ax = \sum_{i=1}^{d} \lambda_i \langle u_i, x \rangle u_i$  with eigenvalues  $\lambda_1 \geq \cdots \geq \lambda_d > 0$  and an orthonormal basis  $u_1, \ldots, u_d$  of V. If we write  $A^{1/2}(x_* - x_0) = \sum_{i=1}^{d} \gamma_i u_i$  for some  $\gamma \in \mathbb{R}^d$ , then

$$egin{aligned} &egin{aligned} &egin{aligned} &egin{aligned} &egin{aligned} &egin{aligned} &A^{1/2}(oldsymbol{x}_*-oldsymbol{x}_k) &iginarrow^2 &=& & \ η_{eta\in\mathcal{Q}_k}\sum_{i=1}^d \gamma_i^2(1-Q(\lambda_i))^2, \end{aligned} \end{aligned}$$

where  $Q_k$  denotes the set of all real polynomials Q of order k with Q(0) = 0.

**Early stopping.** If the set  $\{\lambda_1, \lambda_2, \dots, \lambda_d\}$  consists of the  $k(\mathbf{A}) \leq d$  different numbers  $\lambda_{(1)} > \lambda_{(2)} > \dots > \lambda_{(k(\mathbf{A}))}$ , then

$$Q(t) := 1 - \prod_{j=1}^{k(A)} (1 - t/\lambda_{(j)})$$

defines a polynomial from  $\mathcal{Q}_{k(\mathbf{A})}$  such that Q = 0 on  $\{\lambda_1, \lambda_2, \dots, \lambda_d\}$ . Consequently,

$$\boldsymbol{x}_{k(\boldsymbol{A})} = \boldsymbol{x}_{*}$$

for arbitrary starting point  $x_0$ ! This result should be taken with a pinch of salt, though. Numerical errors may result in more than k(A) iterations, sometimes even more than d. Indeed, one should take precautions to avoid endless loops.

**Tshebyshev polynomials.** The subsequent error bounds use the well-known Tshebyshev polynomials. For  $k \in \mathbb{N}_0$  and  $s \in \mathbb{R}$  let

$$T_k(s) := \begin{cases} (-1)^k \cosh(k \operatorname{arcosh}(-s)) & \text{if } s \leq -1, \\ \cos(k \operatorname{arccos}(s)) & \text{if } s \in [-1, 1], \\ \cosh(k \operatorname{arcosh}(s)) & \text{if } s \geq 1. \end{cases}$$

Then  $T_0(s) = 1, T_1(s) = s$ , and

(6.4) 
$$T_{k+1}(s) = 2sT_k(s) - T_{k-1}(s) \text{ for } k \in \mathbb{N}.$$

The claims about  $T_0$  and  $T_1$  are obvious. The recursion formula (6.4) can be verified as follows: Note first that

$$\cosh(z_1 + z_2) + \cosh(z_1 - z_2) = 2\cosh(z_1)\cosh(z_2)$$
 for  $z_1, z_2 \in \mathbb{C}$ .

Since  $\cos(z) = \cosh(iz)$ , the latter equation is true with  $\cos(\cdot)$  in place of  $\cosh(\cdot)$ . Thus, for  $s \ge 1$ ,

$$T_{k+1}(\pm s) = (\pm 1)^{k+1} \cosh\left(\operatorname{arcosh}(s) + k \operatorname{arcosh}(s)\right)$$
  
=  $2(\pm 1)^{k+1} \cosh\left(\operatorname{arcosh}(s)\right) \cosh\left(k \operatorname{arcosh}(s)\right)$   
-  $(\pm 1)^{k+1} \cosh\left(\operatorname{arcosh}(s) - k \operatorname{arcosh}(s)\right)$   
=  $2(\pm s)(\pm 1)^k \cosh\left(k \operatorname{arcosh}(s)\right) - (\pm 1)^{k-1} \cosh\left((k-1) \operatorname{arcosh}(s)\right)$   
=  $2(\pm s)T_k(\pm s) - T_{k-1}(\pm s).$ 

The case  $s \in [-1, 1]$  can be treated analogously.

The formulae  $T_0(s) = 1$ ,  $T_1(s) = s$  and the recursion formula (6.4) imply that  $T_k(s)$  is equal to  $\sum_{j=0}^k \alpha_{kj} s^j$  with certain coefficients  $\alpha_{kj} \in \mathbb{R}$ , where  $\alpha_{kk} = 2^{k-1}$  for  $k \in \mathbb{N}$ . Moreover, one can show that

$$\begin{aligned} |T_k| &\leq 1 \quad \text{on } [-1,1], \\ |T_k| &> 1 \quad \text{on } \mathbb{R} \setminus [-1,1], \\ T_k \Big( \cos\Big(\frac{2j-1}{k}\pi\Big) \Big) &= 0 \quad \text{for } j = 1, \dots, k, \\ T_k \Big( \cos\Big(\frac{2j}{k}\pi\Big) \Big) &= (-1)^j \quad \text{for } j = 0, 1, \dots, k. \end{aligned}$$

Another useful representation of  $T_k$  results from the known formulae

$$\operatorname{arcosh}(s) = \log(s + \sqrt{s^2 - 1}) = -\log(s - \sqrt{s^2 - 1}) \text{ for } s \ge 1.$$

This implies that

$$T_k(\pm s) = (-1)^k \frac{\left(s + \sqrt{s^2 - 1}\right)^k + \left(s - \sqrt{s^2 - 1}\right)^k}{2} \quad \text{for } s \ge 1.$$

And from this one can deduce that

$$T_k(s) = \sum_{\ell=0}^{\lfloor k/2 \rfloor} \binom{k}{2\ell} s^{k-2\ell} (s^2 - 1)^{\ell}.$$

$$\frac{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{x}_{*}-\boldsymbol{x}_{k})\right\|}{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{x}_{*}-\boldsymbol{x}_{0})\right\|} = \min_{Q \in \mathcal{Q}_{k}} \left(\sum_{i=1}^{d} p_{i}(1-Q(\lambda_{i}))^{2}\right)^{1/2} \leq \min_{Q \in \mathcal{Q}_{k}} \max_{t \in [\lambda_{d},\lambda_{1}]} |1-Q(t)|.$$

As shown later, the minimum is attained for a special polynomial  $Q_k$  such that

$$\max_{t \in [\lambda_d, \lambda_1]} |1 - Q_k(t)| = |1 - Q_k(\lambda_d)| = |1 - Q_k(\lambda_1)|.$$

Hence, if  $oldsymbol{A}^{1/2}(oldsymbol{x}_*-oldsymbol{x}_0)$  fulfils

$$p_i = 0$$
 whenever  $\left|1 - Q_k(\lambda_i)\right| < \max_{t \in [\lambda_d, \lambda_1]} \left|1 - Q_k(t)\right|,$ 

then

$$\frac{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{x}_{*}-\boldsymbol{x}_{k})\right\|}{\left\|\boldsymbol{A}^{1/2}(\boldsymbol{x}_{*}-\boldsymbol{x}_{0})\right\|} \ = \ \max_{t\in[\lambda_{d},\lambda_{1}]} \big|1-Q_{k}(t)\big|.$$

**Theorem 6.6**. For  $k \in \mathbb{N}$  and real numbers 0 < a < b let

$$Q_k(t \mid a, b) := 1 - T_k \Big( \frac{2t - a - b}{b - a} \Big) \Big/ T_k \Big( - \frac{a + b}{b - a} \Big),$$

where  $T_k$  is the Tshebyshev polynomial of order k. Then  $Q_k(\cdot | a, b) \in \mathcal{Q}_k$ , and

$$\max_{t \in [a,b]} \left| 1 - Q_k(t \mid a, b) \right| = \delta_k(a/b)$$

with

$$\delta_k(u) := \frac{2(1-u)^k}{\left(1+\sqrt{u}\right)^{2k} + \left(1-\sqrt{u}\right)^{2k}} \quad \text{for} \quad u \in [0,1].$$

Furthermore, for arbitrary  $Q \in \mathcal{Q}_k$ ,

$$\max_{t \in [a,b]} \left| 1 - Q(t) \right| \ge \delta_k(a/b)$$

with equality if and only if  $Q \equiv Q_k(\cdot \mid a, b)$ .

Figure 6.2 shows for a = 1, b = 4 and k = 2, 3, 6, 7 the optimal polynomial  $Q_k(\cdot | a, b)$ . **Remark 6.7** (Inequalities for  $\delta_k(\cdot)$ ). Obviously,

$$\delta_k(u) < 2 \frac{(1-u)^k}{(1+\sqrt{u})^{2k}} = 2 \left(\frac{1-\sqrt{u}}{1+\sqrt{u}}\right)^k \text{ for } 0 \le u < 1.$$

For  $k \geq 2$  and 0 < u < 1, strict convexity of  $s \mapsto s^k$  on  $(0,\infty)$  implies that

$$\frac{\left(1+\sqrt{u}\right)^{2k}+\left(1-\sqrt{u}\right)^{2k}}{2}=\frac{\left(1+u+2\sqrt{u}\right)^{k}+\left(1+u-2\sqrt{u}\right)^{k}}{2} > (1+u)^{k}.$$

Hence,

$$\delta_k(u) < \delta_1(u)^k \quad \text{for } k \ge 2 \text{ and } 0 < u < 1.$$



Figure 6.2: Some optimal polynomials  $Q_k(\cdot | a, b)$  for a = 1 and b = 4.

**Remark 6.8** (Expansions for  $\delta_k(\cdot)$ ). By virtue of the binomial formula,

$$\frac{\left(1+\sqrt{u}\right)^{2k}+\left(1-\sqrt{u}\right)^{2k}}{2} = \sum_{\ell=0}^{k} \binom{2k}{2\ell} u^{\ell} = \begin{cases} 1+u & \text{for } k=1, \\ 1+k(2k-1)u+O(u^2) & \text{for } k\geq 1. \end{cases}$$

In particular,

$$\delta_1(u) = \frac{1-u}{1+u} = 1-2u+O(u^2),$$

and

$$\delta_1(u)^k = 1 - 2ku + O(u^2),$$
  
$$\delta_k(u) = 1 - 2k^2u + O(u^2).$$

**Remark 6.9** (Comparison of gradient and CG method). Theorem 6.6 and Remark 6.7 imply that k steps of the CG algorithm result in a reduction of the approximation error  $\|A^{1/2}(x_* - x)\|$  by a factor  $\delta_k(\lambda_d/\lambda_1)$  at least. This factor is strictly smaller than the factor  $\delta_1(\lambda_d/\lambda_1)^k$  which could result from k iterations of the gradient method.

162



Figure 6.3: Approximation errors  $\delta_1(\cdot) > \delta_1(\cdot)^k > \delta_k(\cdot)$ .

Figure 6.3 shows the functions  $\delta_1(\cdot) > \delta_1(\cdot)^k > \delta_k(\cdot)$  for k = 2, 3.

In the proof of Theorem 6.6 we use an elementary lemma about real polynomials.

**Lemma 6.10.** For  $k \in \mathbb{N}_0$  let  $\Delta : \mathbb{R} \to \mathbb{R}$  be given by  $\Delta(s) = \sum_{j=0}^k \alpha_j s^j$  with real coefficients  $\alpha_0, \alpha_1, \ldots, \alpha_k$ . Suppose that there exist points  $s_0 < s_1 < \cdots < s_{k+1}$  such that

$$(-1)^i \Delta(s_i) \geq 0$$
 for  $0 \leq i \leq k+1$ .

Then,  $\Delta \equiv 0$ .

Exercise 6.11. Prove Lemma 6.10.

**Proof of Theorem 6.6.** One can easily verify that  $Q_k(\cdot) = Q_k(\cdot | a, b)$  defines a polynomial in  $Q_k$ . Further, the properties of the Tshebyshev polynomials imply that

$$\delta_k := \max_{t \in [a,b]} \left| 1 - Q_k(t) \right| = \left| 1 - Q_k(a) \right| = \frac{|T_k(-1)|}{\left| T_k(-(a+b)/(b-a)) \right|}$$
$$= \frac{1}{\cosh(k \operatorname{arcosh}(1/\delta))}$$

with  $\delta := (b-a)/(a+b) = (1-u)/(1+u)$ . In particular,  $\delta_1 = \delta$ . For general  $k \ge 1$ ,

$$\begin{split} \delta_k &= \frac{2}{\left(\delta^{-1} + \sqrt{\delta^{-2} - 1}\right)^k + \left(\delta^{-1} - \sqrt{\delta^{-2} - 1}\right)^k} \\ &= \frac{2\delta^k}{\left(1 + \sqrt{1 - \delta^2}\right)^k + \left(1 - \sqrt{1 - \delta^2}\right)^k} \\ &= \frac{2(1 - u)^k}{\left(1 + u + \sqrt{(1 + u)^2 - (1 - u)^2}\right)^k + \left(1 + u - \sqrt{(1 + u)^2 - (1 - u)^2}\right)^k} \\ &= \frac{2(1 - u)^k}{\left(1 + u + 2\sqrt{u}\right)^k + \left(1 + u - 2\sqrt{u}\right)^k} \\ &= \frac{2(1 - u)^k}{\left(1 + \sqrt{u}\right)^{2k} + \left(1 - \sqrt{u}\right)^{2k}}. \end{split}$$

It remains to prove optimality and uniqueness of  $Q_k$ . Suppose there exists a polynomial  $Q \in Q_k$  such that

$$\max_{t\in[a,b]} \left|1-Q(t)\right| \leq \delta_k$$

Let  $a = t_1 < t_2 < \cdots < t_{k+1} = b$  be the points  $t_i \in [a, b]$  such that  $Q_k(t_i) = 1 - (-1)^{i-1}\delta_k$ . Then the polynomial  $\Delta := Q_k - Q$  satisfies  $\Delta(0) = 0$  and  $(-1)^i \Delta(t_i) \ge 0$  for  $1 \le i \le k+1$ . Hence, it follows from Lemma 6.10 that  $\Delta \equiv 0$ .

**Example 6.12.** We consider the matrix  $\mathbf{A} = \text{diag}(2^{1-i} : i = 1, ..., d)$  and the vector  $\mathbf{b} = (1, 1, ..., 1)^{\top}$  with d = 8. For the starting point  $\mathbf{x}_0 = \mathbf{0}$ , the CG algorithm yields the following approximations  $\mathbf{x}_k = (x_{k,i})_{i=1}^d$  and errors (all numbers rounded to three digits):

k	$x_{k,1}$	$x_{k,2}$	$x_{k,3}$	$x_{k,4}$	$x_{k,5}$	$x_{k,6}$	$x_{k,7}$	$x_{k,8}$	$\ m{r}_k\ $
0	0	0	0	0	0	0	0	0	2.828
1	4.016	4.016	4.016	4.016	4.016	4.016	4.016	4.016	3.674
2	-0.423	6.287	9.643	11.320	12.159	12.578	12.788	12.893	3.265
3	1.271	-0.984	9.273	17.248	21.947	24.474	25.782	26.448	2.560
4	0.978	2.610	-0.494	14.748	28.789	37.692	42.649	45.258	1.807
5	1.001	1.953	4.859	2.395	24.889	46.736	60.991	69.043	1.128
6	1.000	2.001	3.947	8.852	10.429	43.122	74.498	95.047	0.589
7	1.000	2.000	4.001	7.969	16.492	28.442	74.842	117.072	0.222
8	1.000	2.000	4.000	8.000	16.000	32.000	64.000	128.000	0.000

Note that  $f(\boldsymbol{x}_k)$  and  $\|\boldsymbol{A}^{1/2}(\boldsymbol{x}_* - \boldsymbol{x}_k)\| = \|\boldsymbol{A}^{-1/2}\boldsymbol{r}_k\|$  are non-decreasing in k, but  $\|\boldsymbol{r}_k\|$  can increase every now and then.

**Example 6.13** (Hilbert matrix). A notorious example for a symmetric and positive definite, but ill-conditioned matrix A is the Hilbert matrix

$$\boldsymbol{A} := \left(\frac{1}{i+j-1}\right)_{i,j=1}^{d}.$$

Obviously, A is symmetric. That A is positive definite can be deduced from the representation

$$\frac{1}{i+j-1} = \int_0^1 t^{i-1} t^{j-1} \, dt.$$

$$oldsymbol{h}^{ op}oldsymbol{A}oldsymbol{h} \ = \ \int_0^1 \Bigl(\sum_{i=1}^d h_i t^{i-1}\Bigr)^2 dt \ > \ 0 \quad ext{for} \ oldsymbol{h} \in \mathbb{R}^d \setminus \{oldsymbol{0}\}.$$

However, the ratio  $\lambda_d/\lambda_1$  tends to be very small, even for moderate dimension d. For instance, if d = 8, then  $(\lambda_1, \lambda_d) \approx (1.696, 1.112 \cdot 10^{-10})$ . For  $\mathbf{b} = (1, 1, \dots, 1)^{\top}$ , the equation system  $A\mathbf{x} = \mathbf{b}$  can not be solved with standard methods such as the R procedure gr.solve. But the CG algorithm with starting point  $\mathbf{x}_0 = \mathbf{0}$  yields after k = 19 steps an approximation  $\mathbf{x}_k$  of  $\mathbf{x}_*$  such that  $\|\mathbf{b} - \mathbf{A}\mathbf{x}_k\| < 10^{-8}$ .

**Example 6.14** (Poisson problem). For integers  $d_x, d_y > 2$ , we consider the vector space  $V = \mathbb{R}^{d_x \times d_y}$  with inner product

$$\langle \boldsymbol{x}, \boldsymbol{y} \rangle := \sum_{i=1}^{d_x} \sum_{j=1}^{d_y} x_{ij} y_{ij} = \operatorname{trace}(\boldsymbol{x}^\top \boldsymbol{y}).$$

Occasionally we interpret a matrix  $x \in V$  as a function  $(i, j) \mapsto x(i, j) = x_{ij}$  on the index set  $\mathcal{D} := \{1, 2, \dots, d_x\} \times \{1, 2, \dots, d_y\}.$ 

Here is a discrete version of the classical Poisson problem (heat equation): Let

$$\mathcal{I} := \{2, \dots, d_x - 1\} \times \{2, \dots, d_y - 1\} \text{ and} \\ \mathcal{R} := \{(i, j) \in \mathcal{D} : i \in \{1, d_x\} \text{ oder } j \in \{1, d_y\}\}$$

be the "interior" and the "boundary" of  $\mathcal{D}$ , respectively. For a given matrix  $b \in V$  with

$$\boldsymbol{b} = 0$$
 on  $\mathcal{I}$ 

we look for a matrix  $oldsymbol{x} \in oldsymbol{V}$  such that

$$x = b$$
 on  $\mathcal{R}$ 

and

(6.5) 
$$x_{ij} = \sum_{u,v=-1,0,1} p_{u,v} x_{i+u,j+v} \text{ for all } (i,j) \in \mathcal{I}.$$

Here

$$p_{u,v} := \begin{cases} 3/16 & \text{if } |u| + |v| = 1, \\ 1/16 & \text{if } |u| + |v| = 2, \\ 0 & \text{if } u = v = 0. \end{cases}$$

In other words, we are looking for a solution  $x \in V$  of the equation

Lx = b

with the linear operator L:V
ightarrow V given by

$$(\boldsymbol{L}\boldsymbol{x})_{ij} := \begin{cases} x_{ij} & \text{for } (i,j) \in \mathcal{R}, \\ x_{ij} - \sum_{u,v=-1,0,1} p_{u,v} x_{i+u,j+v} & \text{for } (i,j) \in \mathcal{I}. \end{cases}$$

The operator L is easily implemented, and the computation of Lx requires  $O(d_x d_y)$  steps. However, representating L by a matrix would result in a matrix of size  $d_x d_y \times d_x d_y$ . This can be very large, even for moderate values  $d_x$  and  $d_y$ .

The operator L is nonsingular. Indeed, suppose that Lx = 0 for some  $x \in V$ . Let  $(i, j) \in D$ with  $|x_{ij}| = \max_{i',j'} |x_{i'j'}|$ . Since x = Lx = 0 on  $\mathcal{R}$ , we may assume without loss of generality that  $(i, j) \in \mathcal{I}$ . It follows from (6.5) that  $x_{i+u,j+v} = x_{ij}$  for all  $u, v \in \{-1, 0, 1\}$ . Now one may deduce inductively that x is a constant function (matrix), and since x = Lx = 0 on  $\mathcal{R}$ , this implies that x = 0.

Solution via CG algorithm. Unfortunately, L is not self-adjoint, but the adjoint operator  $L^*$ :  $V \to V$ , given by

$$\langle m{x}, m{L}m{y} 
angle \; = \; \langle m{L}^*m{x}, m{y} 
angle \;\;$$
 for arbitrary  $m{x}, m{y} \in m{V}$  ,

is given by

$$(\boldsymbol{L}^*\boldsymbol{x})_{ij} = x_{ij} - \sum_{u,v=-1,0,1} \mathbf{1}_{[(i+u,j+v)\in\mathcal{I}]} p_{u,v} x_{i+u,j+v}.$$

Thus we consider

$$\widetilde{f}(\boldsymbol{x}) := 2^{-1} \|L\boldsymbol{x} - \boldsymbol{b}\|^2 = 2^{-1} \langle \boldsymbol{x}, L^*L\boldsymbol{x} \rangle - \langle \boldsymbol{x}, L^*\boldsymbol{b} \rangle + 2^{-1} \|\boldsymbol{b}\|^2$$

and apply the CG algorithm with the self-adjoined and positive definite operator  $A := L^*L$  and  $L^*b$  in place of b.

Figure 6.4 depicts a gray scale image of the approximate solution  $x_k$  in case of  $d_x = d_y = 50$  and

$$b_{ij} = \begin{cases} +1 & \text{if } (i,j) \in \mathcal{R} \text{ and } \min\{|i-1|+|j-1|, |i-50|+|j-50|\} \le 9, \\ -1 & \text{if } (i,j) \in \mathcal{R} \text{ and } \min\{|i-1|+|j-50|, |i-50|+|j-1|\} \le 9, \\ 0 & \text{else.} \end{cases}$$

This corresponds to the temperature on a metal plate which is heated at to opposite corners and cooled down at the other two corners. The solution depicted here resulted from the starting point  $x_0 = b$  and the stopping criterion  $||L^*b - L^*Lx_k|| < \epsilon := 10^{-9}$  after k = 761 steps, which is substantially smaller than  $d_x d_y = 2500$ . Figure 6.5 shows  $\log_{10} ||b - Lx_s||$  as a function of  $s \in \{0, 1, \ldots, k\}$ .

Solution via fixed point iteration. The solution  $x_*$  may be interpreted stochastically. let  $(Z_t)_{t\geq 0}$  be a random walk on  $\mathbb{Z} \times \mathbb{Z}$  with independent increments

$$Z_{t+1} - Z_t \sim \sum_{u,v=-1,0,1} p_{u,v} \delta_{(u,v)}$$

and fixed starting value  $Z_0 \in \mathcal{D}$ . Further let T be the stopping time

$$T := \min\{t \in \mathbb{N}_0 : Z_t \in \mathcal{R}\}.$$



Figure 6.4: Approximate solution of the Poisson problem (Example 6.14).

Then,

$$oldsymbol{x}_*(i,j) \;=\; \operatorname{I\!E}ig(oldsymbol{b}(Z_T)\,ig|\,Z_0=(i,j)ig) \quad ext{for}\;(i,j)\in\mathcal{D}.$$

If we define for  $s \in \mathbb{N}_0$  the matrix  $\boldsymbol{y}_s \in \boldsymbol{V}$  via

$$\boldsymbol{y}_s(i,j) := \mathbb{E} \big( \boldsymbol{b}(Z_{\min(T,s)}) \, \big| \, Z_0 = (i,j) \big),$$

then one can show that

$$egin{aligned} oldsymbol{y}_0 &= oldsymbol{b}, \ oldsymbol{y}_{s+1} &= Moldsymbol{y}_s & ext{for } s = 0, 1, 2, \dots, \ \lim_{s o \infty} oldsymbol{y}_s &= oldsymbol{x}_* &= Moldsymbol{x}_*. \end{aligned}$$

Here  $M: \mathbf{V} \to \mathbf{V}$  is the linear operator with

$$(M\boldsymbol{x})_{ij} := \begin{cases} x_{ij} & \text{for } (i,j) \in \mathcal{R}, \\ \sum_{u,v=-1,0,1} p_{u,v} x_{i+u,j+v} & \text{for } (i,j) \in \mathcal{I}. \end{cases}$$

For our particular example, Figure 6.5 shows in addition the log-approximation errors  $\log_{10} \| \boldsymbol{b} - L \boldsymbol{y}_s \|$ ,  $0 \le s \le k = 761$ . For quite some time,  $\| \boldsymbol{b} - L \boldsymbol{y}_s \|$  is smaller than  $\| \boldsymbol{b} - L \boldsymbol{x}_s \|$ , but eventually, as s approaches k,  $\| \boldsymbol{b} - L \boldsymbol{x}_s \|$  is substantially smaller than  $\| \boldsymbol{b} - L \boldsymbol{y}_s \|$ .



Figure 6.5: Log-approximation errors in Poisson problem (Example 6.14): CG algorithm (•) and fixed point iterations (blue line).

#### 6.6 Minimizing a Smooth Convex Function

Let  $\mathcal{X} \subset \mathbb{R}^d$  be an open, convex set, and let  $f : \mathcal{X} \to \mathbb{R}$  be twice continuously differentiable such that its Hessian matrix  $A(x) := D^2 f(x)$  is positive definite for any  $x \in \mathcal{X}$ . Suppose that f has a unique minimizer  $x_* \in \mathcal{X}$ . Now we describe an iterative procedure for the approximation of  $x_*$  which is inspired by the CG algorithm. To this end we write  $r(x) := -\nabla f(x)$ . The idea is to approximate f(x + h) for given  $x \in \mathcal{X}$  and "small"  $h \in \mathbb{R}^d$  by

$$\widetilde{f}(\boldsymbol{x} + \boldsymbol{h} \,|\, \boldsymbol{x}) := f(\boldsymbol{x}) - \boldsymbol{r}(\boldsymbol{x})^{\top} \boldsymbol{h} + 2^{-1} \boldsymbol{h}^{\top} \boldsymbol{A}(\boldsymbol{x}) \boldsymbol{h}.$$

We start with an arbitrary point  $x_0 \in \mathcal{X}$  and set

$$ilde{m{h}}_1 \, := \, rgmin_{m{h}\in {
m span}(m{r}(m{x}_0))} \, \widetilde{f}(m{x}_0 + m{h} \,|\, m{x}_0) \, = \, rac{\|m{r}(m{x}_0)\|^2}{m{r}(m{x}_0)^ op m{A}(m{x}_0)m{r}(m{x}_0)} \, m{r}(m{x}_0).$$

Suppose that for some number  $k \in \mathbb{N}$ , we have already chosen points  $\boldsymbol{x}_0, \boldsymbol{x}_1, \dots, \boldsymbol{x}_{k-1} \in \mathcal{X}$  and directions  $\tilde{\boldsymbol{h}}_1, \tilde{\boldsymbol{h}}_2, \dots, \tilde{\boldsymbol{h}}_k \in \mathbb{R}^d$  such that  $\boldsymbol{r}(\boldsymbol{x}_{k-1})^\top \tilde{\boldsymbol{h}}_k > 0$ . Then we choose

$$oldsymbol{x}_k \ := \ oldsymbol{x}_{k-1} + \lambda_k oldsymbol{ ilde{h}}_k$$

for a certain step size  $\lambda_k = \lambda(\boldsymbol{x}_{k-1}, \tilde{\boldsymbol{h}}_k) \in (0, 1]$  such that  $\boldsymbol{x}_k \in \mathcal{X}$  and  $f(\boldsymbol{x}_k) < f(\boldsymbol{x}_{k-1})$ . Then we define

$$ilde{oldsymbol{h}}_{k+1} := rgmin_{oldsymbol{h}\in ext{span}ig( ilde{oldsymbol{h}}_k, oldsymbol{r}(oldsymbol{x}_k)ig)} ilde{f}(oldsymbol{x}_k+oldsymbol{h}\,|oldsymbol{x}_k).$$

If f is a quadratic function on  $\mathcal{X} = \mathbb{R}^d$ , then this algorithm coincides essentially with the CG algorithm, provided that  $\lambda_k = 1$  for all k. In general, the computation of  $h_{k+1}$  amounts to minimizing a uni- or bivariate quadratic function which is much easier than the computation of a usual Newton step  $A(\mathbf{x}_k)^{-1} \mathbf{r}(\mathbf{x}_k)$ .

## **Chapter 7**

# **Dynamic Programming**

Dynamic programming is a strategy to solve discrete optimization problems. More precisely, one wants to minimize a target function f on a *finite* but huge set  $\mathcal{X}$ . The idea of dynamic programming is to identify a sequence of minimization problems  $(\mathcal{X}_1, f_1), (\mathcal{X}_2, f_2), \ldots, (\mathcal{X}_n, f_n) = (\mathcal{X}, f)$  such that

• minimizing  $f_1$  over  $\mathcal{X}_1$  is easy,

• for  $2 \le k \le n$ , minimizing  $f_k$  over  $\mathcal{X}_k$  is easy, provided that we have already minimized  $f_j$  over  $\mathcal{X}_j$  for  $1 \le j < k$ .

We illustrate this paradigm with three particular examples. For a more thorough treatment we refer to the textbook of Cormen et al. (1990).

#### 7.1 Dykstra's Algorithm

We are given  $n \ge 3$  "locations" which we label with the numbers 1, 2, ..., n. For two locations a, b let  $D(a, b) \in [0, \infty]$  be their "direct distance", where we assume that D(a, a) = 0. Specifically one may think of geographic locations, and D(a, b) could measure the distance, time or costs when travelling from a to b directly. The value  $D(a, b) = \infty$  could indicate that there is no direct route from a to b.

For two locations a and b we are now interested in a path

$$\vec{x} = (x_0, x_1, \dots, x_m)$$

of arbitrary length  $m \ge 1$  consisting of locations  $x_0, x_1, \ldots, x_m$  such that  $x_0 = a$  and  $x_m = b$ with minimal total distance

$$f(\vec{x}) := \sum_{j=1}^{m} D(x_{j-1}, x_j).$$

Since  $D \ge 0$ , it suffices to consider paths of length at most n-1, because constant pieces or loops would never diminish total distance. The assumption that D(c, c) = 0 for all c implies that any path of lenth m < n-1 may be extended by a constant piece without changing its total distance. Hence we consider for m = 1, 2, ..., n - 1 the minimal total distances

$$G_m(a,b) := \min_{\vec{x} \in \mathcal{X}_m(a,b)} f(\vec{x})$$

with  $\mathcal{X}_m(a, b)$  denoting the set of all paths  $\vec{x} = (x_0, x_1, \dots, x_m)$  of length m connecting a and b.

In principle one could determine  $G_{n-1}(a, b)$  by considering all  $n^{n-2}$  paths in  $\mathcal{X}_{n-1}(a, b)$ . Or one goes through all (n-1)! paths  $(x_0, x_1, \ldots, x_{n-1})$  such that  $x_0 = a$  and  $\{x_0, x_1, \ldots, x_{n-1}\} = \{1, 2, \ldots, n\}$  and computes for each  $b \neq a$  the sum  $\sum_{j=1}^{j(b)} D(x_{j-1}, x_j)$  with j(b) denoting the unique index  $j \geq 1$  such that  $x_j = b$ . But even for moderate numbers n this would be way too demanding.

Instead one may solve the problem in  $O(n^3)$  steps with memory  $O(n^2)$  for any fixed location a and simultaneously for all destinations b: Note first that

$$G_1(a,b) = D(a,b).$$

For  $m = 2, 3, \ldots, n-1$  we have the recursion

$$G_m(a,b) = \min_{z=1,\dots,n} (G_{m-1}(a,z) + D(z,b)) \le G_{m-1}(a,b),$$

because the part  $(x_0, \ldots, x_{m-1})$  of an optimal path  $(x_0, x_1, \ldots, x_m) \in \mathcal{X}_m(a, b)$  has to be an optimal path in  $\mathcal{X}_{m-1}(a, x_{m-1})$ . These facts are utilized in Dykstra's algorithm described in Table 7.1: For a given location a we determine the matrix

$$\boldsymbol{G} := \left( G_m(a, y) \right)_{1 \le m \le n-1, 1 \le y \le n} \in [0, \infty]^{(n-1) \times n}$$

inductively for m = 2, 3, ..., n - 1. Thereafter, we determine for an arbitrary b an optimal path  $(x_0, x_1, ..., x_{n-1}) \in \mathcal{X}_{n-1}(a, b)$ . To ease the latter task, the algorithm utilizes a second matrix  $J = (J_m(y))_{1 \le m \le n-1, 1 \le y \le n}$  with  $J_1(y) := a$  and

$$J_m(y) \in \underset{z=1,...,n}{\operatorname{arg\,min}} \left( G_{m-1}(a,z) + D(z,y) \right) \text{ for } 2 \le m \le n-1, 1 \le y \le n.$$

Then  $(x_0, x_1, ..., x_{n-1})$  with  $x_{n-1} := b$  and  $x_{m-1} := x_{J_m(x_m)}$  for m = n - 1, n - 2, ..., 1 yields an optimal path in  $\mathcal{X}_{n-1}(a, b)$ .

If we also require that  $J_m(y) = y$  whenever possible, then the optimal path  $(x_0, x_1, \ldots, x_{n-1}) \in \mathcal{X}_{n-1}(a, b)$  just constructed will have at most one constant piece at the end and can be shortened easily, see Exercise 7.2.

**Remark 7.1**. Dykstra's algorithm is often formulated as an algorithm acting on weighted and directed graphs, see for instance Cormen et al. (1990). Depending on the underlying data structure and distance function, it may be worthwhile to utilize the regions  $E(y) := \{z : D(z, y) < \infty\}$  and write

$$J_{m,y} \leftarrow z_o \in \operatorname*{arg\,min}_{z \in E(y)} (G_{m-1,z} + D(z,y))$$

in Table 7.1.

 $\begin{array}{l} \textbf{Algorithm } \vec{x} \leftarrow \textbf{Dykstra}\big(D(\cdot,\cdot),a,b\big) \\ \textbf{J} \leftarrow (a)_{1 \leq m \leq n-1, 1 \leq y \leq n} \\ \textbf{G} \leftarrow (\infty)_{1 \leq m \leq n-1, 1 \leq y \leq n} \\ (G_{1,y})_{1 \leq y \leq n} \leftarrow \big(D(a,y)\big)_{1 \leq y \leq n} \\ \textbf{for } m \leftarrow 2 \textbf{ to } (n-1) \textbf{ do} \\ \textbf{for } y \in \{1, \dots, n\} \textbf{ do} \\ J_{m,y} \leftarrow z_o \in \operatorname*{arg\,min}_{1 \leq z \leq n} \big(G_{m-1,z} + D(z,y)\big) \\ G_{m,y} \leftarrow G_{m-1,z_o} + D(z_o,y) \\ \textbf{end for} \\ \textbf{end for} \\ \vec{x} = (x_0, \dots, x_{n-2}, x_{n-1}) \leftarrow (*, \dots, *, b) \\ \textbf{for } m \leftarrow (n-1) \textbf{ downto 1 do} \\ x_{m-1} \leftarrow J_{m,x_m} \\ \textbf{end for} \end{array}$ 

Table 7.1: Dykstra's Algorithm.

**Exercise 7.2.** Suppose that for arbitrary  $y \in \{1, ..., n\}$  and  $m \in \{2, ..., n-1\}$ ,  $J_m(y) = y$  whenever possible. Show that Dykstra's algorithm produces an optimal path  $(x_0, x_1, ..., x_{n-1})$  in  $\mathcal{X}_{n-1}(a, b)$  such that for some  $M \in \{0, 1, ..., n-1\}$ , the points  $x_j, 0 \le j \le M$  are distinct while  $x_M = \cdots = x_{n-1} = b$ .

**Exercise 7.3** (Negative costs). If one considers the planning of a bike route and if D(a, b) measures strain to get from a to b, it could be reasonable to consider negative numbers D(a, b), for instance, if the direct road from a to b goes mostly downhill. Formulate conditions on a cost function  $D : \{1, \ldots, n\} \times \{1, \ldots, n\} \rightarrow (-\infty, \infty]$  such that Dykstra's algorithm works without modifications.

**Exercise 7.4** (A small world). Consider *n* people which we label with 1, 2, ..., n. Let  $K \in \{0, 1\}^{n \times n}$  be a matrix such that  $K_{a,b} = 1$  indicates that person *a* knows person *b*. Now let  $U \in \{0, 1\}^{n \times n}$  have entry  $U_{a,b} = 1$  if and only if there is a tuple  $(x_j)_{j=0}^m$  of  $m + 1 \le n$  persons such that  $x_0 = a, x_m = b$  and  $K_{x_{j-1},x_j} = 1$  for  $1 \le j \le m$ . How could one determine this matrix U efficiently from K?

#### 7.2 Alignment of Sequences

We are given two sequences  $x = (x_1, x_2, ..., x_m)$  and  $y = (y_1, y_2, ..., y_n)$  of "letters"  $x_i$  and  $y_j$  in an alphabet  $\mathcal{A}$ . For instance, the sequences x and y could represent amino acids of two proteins, i.e.  $\mathcal{A}$  consists of the 20 possible amino acids. Alternatively, x and y could be two DNA sequences, represented by letters in the alphabet  $\mathcal{A}$  of 4 possible nucleic acids.

To judge how similar these sequences x and y are, we want to determine an optimal 'alignment' of them. That means, we are looking for a matrix

$$\boldsymbol{H} = \begin{bmatrix} H_{1,1} & H_{1,2} & \cdots & H_{1,N} \\ H_{2,1} & H_{2,2} & \cdots & H_{2,N} \end{bmatrix}$$

with entries  $H_{i,j} \in \overline{A} := A \cup \{\diamond\}$ , where  $\diamond \notin A$  represents a blank (space), such that the two rows  $(H_{1,j})_{j=1}^N$  and  $(H_{2,j})_{j=1}^N$  contain the entries of x and y, respectively, in their original order, augmented by several blanks  $\diamond$  at arbitrary places. The number N of columns is arbitrary but will certainly satisfy  $N \ge \max(m, n)$ . Our goal is to minimize

$$f(\boldsymbol{H}) := \sum_{j=1}^{N} c(H_{1,j}, H_{2,j})$$

over all alignments  $\boldsymbol{H}$  of  $\boldsymbol{x}$  and  $\boldsymbol{y}$ , where  $c : \bar{\mathcal{A}} \times \bar{\mathcal{A}} \to [0, \infty)$  is a certain cost function satisfying  $c(\diamond, \diamond) = 0$ ,  $c(a, a) \leq 0$  for all  $a \in \mathcal{A}$ , and c(a, b) > 0 for  $a, b \in \bar{\mathcal{A}}$  with  $a \neq b$ . Obviously it suffices to consider alignments  $\boldsymbol{H}$  with  $N \leq m + n$  columns, because columns with two blanks  $\diamond$  could be deleted without increasing  $f(\boldsymbol{H})$ .

To construct an optimal alignment  $\boldsymbol{H}$  via dynamic programming we consider the subsequences  $\boldsymbol{x}^{(k)} := (x_1, x_2, \ldots, x_k)$  and  $\boldsymbol{y}^{(\ell)} := (y_1, y_2, \ldots, y_\ell)$ , including the empty sequences  $\boldsymbol{x}^{(0)}$  and  $\boldsymbol{y}^{(0)}$ . Now one can easily verify that an optimal alignment  $\boldsymbol{H} \in \bar{\mathcal{A}}^{2 \times N}$  of  $\boldsymbol{x}^{(k)}$  and  $\boldsymbol{y}^{(\ell)}$  with  $k, \ell \geq 1$  needs to satisfy one of the following three conditions:

- (i)  $H_{1,N} = x_k, H_{2,N} = y_\ell$ , and  $(H_{i,j})_{1 \le i \le 2, 1 \le j < N}$  is an optimal alignment of  $x^{(k-1)}$  and  $y^{(\ell-1)}$ ;
- (ii)  $H_{1,N} = x_k, H_{2,N} = \diamond$ , and  $(H_{i,j})_{1 \le i \le 2, 1 \le j < N}$  is an optimal alignment of  $\boldsymbol{x}^{(k-1)}$  and  $\boldsymbol{y}^{(\ell)}$ ;
- (iii)  $H_{1,N} = \diamond, H_{2,N} = y_{\ell},$ and  $(H_{i,j})_{1 \le i \le 2, 1 \le j < N}$  is an optimal alignment of  $\boldsymbol{x}^{(k)}$  and  $\boldsymbol{y}^{(\ell-1)}$ .

Hence we construct a matrix  $G = (G_{k,\ell})_{0 \le k \le m, 0 \le \ell \le n}$  as follows:

$$G_{0,0} := 0,$$
  

$$G_{k,0} := \sum_{i=1}^{k} c(x_i, \diamond) \quad \text{for } 1 \le k \le m,$$
  

$$G_{0,\ell} := \sum_{j=1}^{\ell} c(\diamond, y_j) \quad \text{for } 1 \le \ell \le n,$$

and for  $1 \le k \le m, 1 \le \ell \le n$  let

$$G_{k,\ell} := \min \left\{ \sum_{j=1}^{N} c(H_{1,j}, H_{2,j}) : \boldsymbol{H} \text{ Alignment of } \boldsymbol{x}^{(k)} \text{ and } \boldsymbol{y}^{(\ell)} \right\}$$
  
= min $\left( G_{k-1,\ell-1} + c(x_k, y_\ell), G_{k-1,\ell} + c(x_k, \diamond), G_{k,\ell-1} + c(\diamond, y_\ell) \right)$ 

This matrix may be computed inductively in O(mn) steps, and one obtains an optimal alignment by backtracking the matrix G as described in Table 7.2.

**Remark 7.5.** Alignment algorithms as the one described here are a standard tool in many genome data bases. Instead of minimizing a cost function, these programs often maximize a similarity function, i.e. c(a, b) measures similarity of two letters  $a, b \in \overline{A}$ . In this case one has to use the assignment

$$G_{k,\ell} \leftarrow \max \left( G_{k-1,\ell-1} + c(x_k, y_\ell), G_{k-1,\ell} + c(x_k, \diamond), G_{k,\ell-1} + c(\diamond, y_\ell) \right)$$

in Table 7.2.

**Exercise 7.6** (Longest monotone subsequences). Let  $x = (x_j)_{j=1}^n$  be a given vector in  $\mathbb{R}^n$ . We are looking for a subvector  $(x_{j(i)})_{i=1}^m$  satisfying

$$j(1) < j(2) < \dots < j(m)$$
 and  $x_{j(1)} \le x_{j(2)} \le \dots \le x_{j(m)}$ 

such that m is as large as possible. Design and implement a dynamic program for this task.

Algorithmus  $oldsymbol{H} \leftarrow \operatorname{Align}(oldsymbol{x},oldsymbol{y},c(\cdot,\cdot))$  $(m, n) \leftarrow (\text{length}(\boldsymbol{x}), \text{length}(\boldsymbol{y}))$  $oldsymbol{G} \leftarrow ig(\inftyig)_{0 \leq k \leq m, 0 \leq \ell \leq n}$  $G_{0,0} \leftarrow 0$  $\left(G_{k,0}\right)_{1\leq k\leq m} \leftarrow \left(\sum_{i=1}^{k} c(x_i,\diamond)\right)_{1\leq k\leq m}$  $\left( G_{0,\ell} \right)_{1 \le \ell \le n} \leftarrow \left( \sum_{j=1}^{\ell} c(\diamond, y_j) \right)_{1 < \ell \le n}^{--}$ for  $k \leftarrow 1$  to m do for  $\ell \leftarrow 1$  to n do  $G_{k,\ell} \leftarrow \min(G_{k-1,\ell-1} + c(x_k, y_\ell), G_{k-1,\ell} + c(x_k, \diamond), G_{k,\ell-1} + c(\diamond, y_\ell))$ end for end for  $\pmb{H} \leftarrow \big(\diamond\big)_{1 \leq i \leq 2, 1 \leq j \leq m+n}$  $(k, \ell) \leftarrow (m, n)$  $T \leftarrow m + n$ while k > 0 and  $\ell > 0$  do if  $G_{k,\ell} = G_{k-1,\ell-1} + c(x_k, y_\ell)$  then  $(H_{1,T}, H_{2,T}) \leftarrow (x_k, y_\ell)$  $(k, \ell) \leftarrow (k - 1, \ell - 1)$ else if  $G_{k,\ell} = G_{k-1,\ell} + c(x_k,\diamond)$  then  $(H_{1,T}, H_{2,T}) \leftarrow (x_k, \diamond)$  $k \leftarrow k - 1$ else  $(H_{1,T}, H_{2,T}) \leftarrow (\diamond, y_{\ell})$  $\ell \leftarrow \ell - 1$ end if  $T \leftarrow T-1$ end while if k > 0 then  $(H_{1,T-k+1},\ldots,H_{1,T}) \leftarrow (x_1,\ldots,x_k)$  $T \leftarrow T - k$ else if  $\ell > 0$  then  $(H_{2,T-\ell+1},\ldots,H_{2,T}) \leftarrow (y_1,\ldots,y_\ell)$  $T \leftarrow T - \ell$ end if  $\boldsymbol{H} \leftarrow \left(H_{i,j}\right)_{1 \le i \le 2, T \le j \le m+n}$ 

Table 7.2: Alignment of two sequences.

### 7.3 Bimonotone Regression

Our third example combines ideas from previous chapters, notably convex polyhedral cones and the PAVA, with dynamic programming to solve a particular regression problem. This material is taken from Beran and Dümbgen [2].

For integers m, n > 1, let  $\boldsymbol{y} \in \mathbb{R}^{m \times n}$  and  $\boldsymbol{w} \in (0, \infty)^{m \times n}$  be given matrices, and consider the function  $f : \mathbb{R}^{m \times n} \to \mathbb{R}$  given by

$$f(x) := \sum_{i,j} w_{ij} (y_{ij} - x_{ij})^2,$$

where  $\sum_{i,j} \cdots$  is shorthand notation for  $\sum_{i=1}^{m} \sum_{j=1}^{n} \cdots$ . As in previous chapters, we identify  $\mathbb{R}^{m \times n}$  with the Euclidean space  $\mathbb{R}^d$ , where d = mn, so the standard inner product is given by  $\langle \boldsymbol{x}, \tilde{\boldsymbol{x}} \rangle := \operatorname{trace}(\boldsymbol{x}^{\top} \tilde{\boldsymbol{x}})$ , and  $\|\boldsymbol{x}\| = \operatorname{trace}(\boldsymbol{x}^{\top} \boldsymbol{x})^{1/2} = \left(\sum_{i,j} x_{ij}^2\right)^{1/2}$  is the Frobenius norm.

Obviously, f is strictly convex and coercive, so there exists a unique minimizer  $x^* \in \mathcal{K}$  of f over any closed, convex subset  $\mathcal{K}$  of  $\mathbb{R}^{r \times s}$ . Note also that

$$\nabla f(\boldsymbol{x}) = 2 \left( w_{ij}(x_{ij} - y_{ij}) \right)_{i \le r, j \le s}$$

in the sense that

$$f(oldsymbol{x}+oldsymbol{v}) \ = \ f(oldsymbol{x}) + \langle 
abla f(oldsymbol{x}), oldsymbol{v} 
angle + \sum_{i,j} w_{ij} v_{ij}^2.$$

Thus a point  $x \in \mathcal{K}$  is equal to  $x^*$  if and only if

(7.1) 
$$\langle \nabla f(\boldsymbol{x}), \tilde{\boldsymbol{x}} - \boldsymbol{x} \rangle \geq 0 \quad \text{for all } \tilde{\boldsymbol{x}} \in \mathcal{K}.$$

Now we consider the special set  $\mathcal{K}$  of all matrices  $x \in \mathbb{R}^{r \times s}$  which are bimonotone in the sense that

$$x_{ij} \le x_{i+1,j}$$
 for  $1 \le i < r$  and  $1 \le j \le s$ ,  
 $x_{ij} \le x_{i,j+1}$  for  $1 \le i \le r$  and  $1 \le j < s$ .

Note that  $\mathcal{K}$  is a closed, concex cone (and a convex polyhedron). Consequently, characterization (7.1) of  $x^*$  is equivalent to the following two conditions:

(7.2) 
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{x} \rangle = 0,$$

(7.3) 
$$\langle \nabla f(\boldsymbol{x}), \tilde{\boldsymbol{x}} \rangle \geq 0 \quad \text{for all } \tilde{\boldsymbol{x}} \in \mathcal{K}.$$

Concerning (7.3), we may apply Exercise 2.37 to show that  $\mathcal{K}$  consists of all matrices

$$x = \sum_{e \in \mathcal{E}} \lambda_e e,$$

where

$$\mathcal{E} := \mathcal{K} \cap \{0, 1\}^{r \times s},$$

and  $\lambda_e \geq 0$  whenever  $e \neq 1 := (1)_{i \leq r, j \leq s}$ . This representation of  $\mathcal{K}$  implies that (7.3) is equivalent to the following two conditions:

(7.4) 
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{1} \rangle = 0,$$

(7.5) 
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{e} \rangle \geq 0 \text{ for all } \boldsymbol{e} \in \mathcal{E}.$$

Now we are ready to describe an explicit active set algorithm for the minimization of f over  $\mathcal{K}$ . It consists of alternating between the following two basic procedures finitely many times.

#### **Basic procedure 1: Finding a locally optimal matrix.** For $x \in \mathcal{K}$ consider the set

$$\mathcal{K}(\boldsymbol{x}) := \{ \tilde{\boldsymbol{x}} \in \mathbb{R}^{m \times n} : \tilde{x}_{ij} \leq \tilde{x}_{k\ell} \text{ whenever } x_{ij} \leq x_{k\ell} \}.$$

One can easily verify that  $\mathcal{K}(\boldsymbol{x})$  is a subcone of  $\mathcal{K}$ . More importantly, if  $\sigma(1), \sigma(2), \ldots, \sigma(mn)$  is a list of all index pairs  $(i, j) \in \{1, \ldots, m\} \times \{1, \ldots, n\}$  such that  $x_{\sigma(1)} \leq x_{\sigma(2)} \leq \cdots \leq x_{\sigma(mn)}$ , then minimizing  $f(\cdot)$  over  $\mathcal{K}(\boldsymbol{x})$  is equivalent to minimizing

$$\sum_{i=1}^{mn} \boldsymbol{w}_{\sigma(i)} (y_{\sigma(i)} - z_i)^2$$

over all vectors  $z \in \mathbb{R}^{mn}$  such that  $z_i \leq z_j$  whenever  $i \leq j$  and  $z_i = z_j$  whenever  $x_{\sigma(i)} = x_{\sigma(j)}$ . This can be done via a suitable modification of the pool-adjacent-violators algorithm. Finally, the resulting minimizer z corresponds to a matrix  $x^{\text{new}} \in \mathcal{K}(x)$  with components  $x_{\sigma(i)}^{\text{new}} = z_i$  for  $1 \leq i \leq mn$ . This matrix satisfies  $f(x^{\text{new}}) \leq f(x)$  with equality if and only if  $x^{\text{new}} = x$ . Finally, we replace x with  $x^{\text{new}}$  and note that the new matrix x is locally optimal in the sense that  $f(\tilde{x}) \geq f(x)$  for all  $\tilde{x} \in \mathcal{K}(x)$ .

Basic procedure 2: Checking for global optimality and moving on, if necessary. Suppose that  $x \in \mathcal{K}$  is locally optimal. Applying (7.2) and (7.4) with  $\mathcal{K}(x)$  in place of  $\mathcal{K}$  shows that x satisfies these two conditions. Consequently, x is equal to the global minimizer  $x^*$  if and only if it satisfies (7.5). To this end, we determine a matrix

$$oldsymbol{e}(oldsymbol{x}) \ \in \ rgmin_{oldsymbol{e}\in\mathcal{E}} lpha \langle 
abla f(oldsymbol{x}), oldsymbol{e} 
angle.$$

If  $\langle \nabla f(\boldsymbol{x}), \boldsymbol{e} \rangle = 0$ , we know that  $\boldsymbol{x} = \boldsymbol{x}^*$ . Otherwise, we replace  $\boldsymbol{x}$  with

$$\boldsymbol{x} + t(\boldsymbol{x})\boldsymbol{e}(\boldsymbol{x})$$

with

$$t(\boldsymbol{x}) := \operatorname*{arg\,min}_{t \geq 0} \left. f(\boldsymbol{x} + t \boldsymbol{e}(\boldsymbol{x})) \right. = \left. -2^{-1} \left\langle 
abla f(\boldsymbol{x}), \boldsymbol{e}(\boldsymbol{x}) \right\rangle \right/ \sum_{i,j} e_{ij}(\boldsymbol{x})$$

The new vector x has a strictly smaller value f(x) than its predecessor. Then we run basic procedure 1 with this new vector x.

In high-dimensional settings, one could also stop earlier when  $-\langle \nabla f(\boldsymbol{x}), \boldsymbol{e}(\boldsymbol{x}) \rangle$  is smaller than a small constant  $\delta > 0$ .

**Minimizing**  $\langle \boldsymbol{a}, \cdot \rangle$  over  $\mathcal{E}$ . With  $\boldsymbol{a} := \nabla f(\boldsymbol{x})$ , basic procedure 2 involves the minimization of  $\langle \boldsymbol{a}, \boldsymbol{e} \rangle$  with respect to  $\boldsymbol{e} \in \mathcal{E}$ . Note that a matrix  $\boldsymbol{e} \in \{0, 1\}^{m \times n}$  belongs to  $\mathcal{K}$  if and only if

$$e_{ij} = 1_{[j>f(i)]}$$

for some non-increasing function  $f : \{1, ..., m\} \to \{0, 1, ..., n\}$ , and a simple combinatorial argument shows that

$$\#\mathcal{E} = \binom{m+n}{m}.$$

Since this grows exponentially in  $\min(m, n)$ , computing  $\langle a, e \rangle$  for all  $e \in \mathcal{E}$  is not feasible in general. Fortunately, the optimization problem can be solved via dynamic programming as follows.

Let  $H = (H_{k,\ell})_{1 \le k \le m+1, 1 \le \ell \le n+1}$  be given by

$$\begin{split} H_{m+1,\cdot} &:= 0, \\ H_{k,n+1} &:= \min \left\{ \sum_{i=k}^{m} \sum_{j=1}^{n} a_{ij} e_{ij} : \boldsymbol{e} \in \mathcal{E} \right\} \quad \text{for } 1 \le k \le m, \\ H_{k,\ell} &:= \min \left\{ \sum_{i=k}^{m} \sum_{j=1}^{n} a_{ij} e_{ij} : \boldsymbol{e} \in \mathcal{E}, e_{k\ell} = 1 \right\} \quad \text{for } 1 \le k \le m, 1 \le \ell \le n, \end{split}$$

so

$$H_{1,n+1} = \min_{\boldsymbol{e}\in\mathcal{E}} \langle \boldsymbol{a}, \boldsymbol{e} \rangle.$$

To compute H via dynamic programming, we use the auxiliary matrix  $B = (B_{i,\ell})_{1 \le i \le m, 1 \le \ell \le n+1}$ with components  $B_{i,n+1} = 0$  and

$$B_{i,\ell} := \sum_{j=\ell}^n a_{ij} \text{ for } 1 \le i \le m, 1 \le \ell \le n.$$

Then for  $1 \le k \le m$  and  $2 \le \ell \le n+1$ ,

$$H_{k,1} = B_{k,1} + H_{k+1,1},$$
  
$$H_{k,\ell} = \min(H_{k,\ell-1}, B_{k,\ell} + H_{k,\ell})$$

Consequently, **B** and **H** can be computed in time O(mn). A minimizer  $e \in \mathcal{E}$  of  $\langle a, e \rangle$  can be

determined by backtracking:

```
\begin{array}{l} \boldsymbol{e} \leftarrow \boldsymbol{0} \\ \boldsymbol{k} \leftarrow 1 \\ \ell \leftarrow n \\ \textbf{while } \boldsymbol{k} \leq m \textbf{ and } \ell \geq 1 \textbf{ do} \\ \textbf{ if } H_{k,\ell} = H_{k,\ell+1} \textbf{ then} \\ (e_{i\ell})_{k \leq i \leq m} \leftarrow 1 \\ \ell \leftarrow \ell - 1 \\ \textbf{ else} \\ \boldsymbol{k} \leftarrow k+1 \\ \textbf{ end if} \\ \textbf{ end while} \end{array}
```

Numerical example. Figure 7.1 shows a numerical example for bimonotone regression with m = n = 100. The upper panel shows the data matrix y on a colour scale from light yellow to dark purple for  $-5 \le y_{ij} \le 5$ . The lower panel depicts the matrix  $x^*$  on a colour scale from light yellow to dark purple for  $-2.2 \le x_{ij} \le 2.7$ .


Figure 7.1: Example for bimonotone regression.