

Chapter 11

Numerical Solution of Generalized Minimax Problems



Ladislav Lukšan, Ctirad Matonoha, and Jan Vlček

Abstract This contribution contains the description and investigation of three numerical methods for solving generalized minimax problems. These problems consists in the minimization of nonsmooth functions which are compositions of special smooth convex functions with maxima of smooth functions. The most important functions of this type are the sums of maxima of smooth functions. Section 11.2 is devoted to primal interior point methods which use solutions of nonlinear equations for obtaining minimax vectors. Section 11.3 contains investigation of smoothing methods, based on using exponential smoothing terms. Section 11.4 contains short description of primal-dual interior point methods based on transformation of generalized minimax problems to general nonlinear programming problems. Finally the last section contains results of numerical experiments.

11.1 Generalized Minimax Problems

In many practical problems we need to minimize functions that contain absolute values or pointwise maxima of smooth functions. Such functions are nonsmooth but they often have a special structure enabling the use of special methods that are more efficient than methods for minimization of general nonsmooth functions. The classical minimax problem, where $F(\mathbf{x}) = \max_{1 \leq k \leq m} f_k(\mathbf{x})$, or problems where the function to be minimized is a nonsmooth norm, e.g. $F(\mathbf{x}) = \|f(\mathbf{x})\|_\infty$, $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_\infty$, $F(\mathbf{x}) = \|f(\mathbf{x})\|_1$, $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_1$ with $f(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_m(\mathbf{x})]^T$ and $f_+(\mathbf{x}) = [\max\{f_1(\mathbf{x}), 0\}, \dots, \max\{f_m(\mathbf{x}), 0\}]^T$, are typical examples. Such functions can be considered as special cases of more general functions, so it is possible to formulate more general theories and construct more general numerical methods. One possibility for generalization of the classical

L. Lukšan (✉) · C. Matonoha · J. Vlček
Institute of Computer Science, The Czech Academy of Sciences, Prague, Czech Republic
e-mail: luksan@cs.cas.cz; matonoha@cs.cas.cz; vlcek@cs.cas.cz

minimax problem consists in the use of the function

$$F(\mathbf{x}) = \max_{1 \leq k \leq \bar{k}} \mathbf{p}_k^T \mathbf{f}(\mathbf{x}), \quad (11.1)$$

where $\mathbf{p}_k \in \mathbb{R}^m$, $1 \leq k \leq \bar{k}$, and $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a smooth mapping. This function is a special case of composite nonsmooth functions of the form $F(\mathbf{x}) = f_0(\mathbf{x}) + \max_{1 \leq k \leq \bar{k}} (\mathbf{p}_k^T \mathbf{f}(\mathbf{x}) + b_k)$, where $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable function [8, Section 14.1].

Remark 11.1 We can express all above mentioned minimax problems and nonsmooth norms in form (11.1).

- (a) Setting $\mathbf{p}_k = \mathbf{e}_k$, where \mathbf{e}_k is the k -th column of an identity matrix and $\bar{k} = m$, we obtain $F(\mathbf{x}) = \max_{1 \leq k \leq m} f_k(\mathbf{x})$ (the classical minimax).
- (b) Setting $\mathbf{p}_k = \mathbf{e}_k$, $\mathbf{p}_{m+k} = -\mathbf{e}_k$ and $\bar{k} = 2m$, we obtain $F(\mathbf{x}) = \max_{1 \leq k \leq m} \max\{f_k(\mathbf{x}), -f_k(\mathbf{x})\} = \|\mathbf{f}(\mathbf{x})\|_\infty$.
- (c) Setting $\mathbf{p}_k = \mathbf{e}_k$, $\mathbf{p}_{m+1} = \mathbf{0}$ and $\bar{k} = m + 1$, we obtain $F(\mathbf{x}) = \max\{\max_{1 \leq k \leq m} f_k(\mathbf{x}), 0\} = \|\mathbf{f}(\mathbf{x})_+\|_\infty$.
- (d) If $\bar{k} = 2^m$ and \mathbf{p}_k , $1 \leq k \leq 2^m$, are mutually different vectors whose elements are either 1 or -1 , we obtain $F(\mathbf{x}) = \sum_{k=1}^m \max\{f_k(\mathbf{x}), -f_k(\mathbf{x})\} = \|\mathbf{f}(\mathbf{x})\|_1$.
- (e) If $\bar{k} = 2^m$ and \mathbf{p}_k , $1 \leq k \leq 2^m$, are mutually different vectors whose elements are either 1 or 0, we obtain $F(\mathbf{x}) = \sum_{k=1}^m \max\{f_k(\mathbf{x}), 0\} = \|\mathbf{f}_+(\mathbf{x})\|_1$.

Remark 11.2 Since the mapping $\mathbf{f}(\mathbf{x})$ is continuously differentiable, the function (11.1) is Lipschitz. Thus, if the point $\mathbf{x} \in \mathbb{R}^n$ is a local minimum of $F(\mathbf{x})$, then $\mathbf{0} \in \partial F(\mathbf{x})$ [25, Theorem 3.2.5] holds. According to [25, Theorem 3.2.13], one has

$$\partial F(\mathbf{x}) = (\nabla \mathbf{f}(\mathbf{x}))^T \text{conv} \{ \mathbf{p}_k : k \in \bar{\mathcal{I}}(\mathbf{x}) \},$$

where $\bar{\mathcal{I}}(\mathbf{x}) = \{k \in \{1, \dots, \bar{k}\} : \mathbf{p}_k^T \mathbf{f}(\mathbf{x}) = F(\mathbf{x})\}$. Thus, if the point $\mathbf{x} \in \mathbb{R}^n$ is a local minimum of $F(\mathbf{x})$, then multipliers $\lambda_k \geq 0$, $1 \leq k \leq \bar{k}$, exist, such that $\lambda_k (\mathbf{p}_k^T \mathbf{f}(\mathbf{x}) - F(\mathbf{x})) = 0$, $1 \leq k \leq \bar{k}$,

$$\sum_{k=1}^{\bar{k}} \lambda_k = 1 \quad \text{and} \quad \sum_{k=1}^{\bar{k}} \lambda_k J^T(\mathbf{x}) \mathbf{p}_k = \mathbf{0},$$

where $J(\mathbf{x})$ is a Jacobian matrix of the mapping $\mathbf{f}(\mathbf{x})$.

Remark 11.3 It is clear that a minimum of function (11.1) is a solution of a nonlinear programming problem consisting in minimization of a function $\tilde{F} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$, where $\tilde{F}(\mathbf{x}, z) = z$, on the set

$$C = \{(\mathbf{x}, z) \in \mathbb{R}^{n+1} : \mathbf{p}_k^T \mathbf{f}(\mathbf{x}) \leq z, 1 \leq k \leq \bar{k}\}.$$

Denoting $c_k(\mathbf{x}, z) = \mathbf{p}_k^T \mathbf{f}(\mathbf{x}) - z$ and $\mathbf{a}_k = \nabla c_k(\mathbf{x}, z)$, $1 \leq k \leq \bar{k}$, we obtain $\mathbf{a}_k = [\mathbf{p}_k^T J(\mathbf{x}), -1]^T$ and $\mathbf{g} = \nabla \tilde{F}(\mathbf{x}, z) = [\mathbf{0}^T, 1]^T$, so the necessary KKT conditions can be written in the form

$$\begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix} + \sum_{k=1}^{\bar{k}} \begin{bmatrix} J^T(\mathbf{x}) \mathbf{p}_k \\ -1 \end{bmatrix} \lambda_k = \mathbf{0},$$

$\lambda_k(\mathbf{p}_k^T \mathbf{f}(\mathbf{x}) - z) = 0$, where $\lambda_k \geq 0$, $1 \leq k \leq \bar{k}$, are the Lagrange multipliers and $z = F(\mathbf{x})$. Thus, we obtain the same necessary conditions for an extremum as in Remark 11.2.

From the examples given in Remark 11.1 it follows that composite nondifferentiable functions are not suitable for representation of the functions $F(\mathbf{x}) = \|f(\mathbf{x})\|_1$ and $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_1$ because in this case the expression on the right-hand side of (11.1) contains 2^m elements with vectors \mathbf{p}_k , $1 \leq k \leq 2^m$. In the subsequent considerations, we will choose a somewhat different approach. We will consider generalized minimax functions established in [5] and [23].

Definition 11.1 We say that $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is a generalized minimax function if

$$F(\mathbf{x}) = h(F_1(\mathbf{x}), \dots, F_m(\mathbf{x})), \quad F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x}), \quad 1 \leq k \leq m, \tag{11.2}$$

where $h : \mathbb{R}^m \rightarrow \mathbb{R}$ and $f_{kl} : \mathbb{R}^n \rightarrow \mathbb{R}$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, are smooth functions satisfying the following assumptions.

Assumption 11.1 Functions f_{kl} , $1 \leq k \leq m$, $1 \leq l \leq m_k$, are bounded from below on \mathbb{R}^n , so that there exists a constant $\underline{F} \in \mathbb{R}$ such that $f_{kl}(\mathbf{x}) \geq \underline{F}$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, for all $\mathbf{x} \in \mathbb{R}^n$.

Assumption 11.2 Functions F_k , $1 \leq k \leq m$, are bounded from below on \mathbb{R}^n , so that there exist constants $\underline{F}_k \in \mathbb{R}$ such that $F_k(\mathbf{x}) \geq \underline{F}_k$, $1 \leq k \leq m$, for all $\mathbf{x} \in \mathbb{R}^n$.

Assumption 11.3 The function h is twice continuously differentiable and convex satisfying

$$0 < \underline{h}_k \leq \frac{\partial}{\partial z_k} h(\mathbf{z}) \leq \bar{h}_k, \quad 1 \leq k \leq m, \tag{11.3}$$

for every $\mathbf{z} \in \mathcal{Z} = \{\mathbf{z} \in \mathbb{R}^m : z_k \geq \underline{F}_k, 1 \leq k \leq m\}$ (vector $\mathbf{z} \in \mathbb{R}^m$ is called the minimax vector).

Assumption 11.4 Functions $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, are twice continuously differentiable on the convex hull of the level set

$$\mathcal{D}_F(\bar{F}) = \{\mathbf{x} \in \mathbb{R}^n : F_k(\mathbf{x}) \leq \bar{F}, 1 \leq k \leq m\}$$

for a sufficiently large upper bound \bar{F} and subsequently, constants \bar{g} and \bar{G} exist such that $\|\mathbf{g}_{kl}(\mathbf{x})\| \leq \bar{g}$ and $\|G_{kl}(\mathbf{x})\| \leq \bar{G}$ for all $1 \leq k \leq m$, $1 \leq l \leq m_k$, and $\mathbf{x} \in \text{conv } \mathcal{D}_F(\bar{F})$, where $\mathbf{g}_{kl}(\mathbf{x}) = \nabla f_{kl}(\mathbf{x})$ and $G_{kl}(\mathbf{x}) = \nabla^2 f_{kl}(\mathbf{x})$.

Remark 11.4 The conditions imposed on the function $h(\mathbf{z})$ are relatively strong but many important nonsmooth functions satisfy them.

- (1) Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be an identity mapping, so $h(z) = z$ and $h'(z) = 1 > 0$. Then setting $\bar{k} = 1$, $m_1 = \bar{l}$ and

$$F(\mathbf{x}) = h(F_1(\mathbf{x})) = F_1(\mathbf{x}) = \max_{1 \leq l \leq \bar{l}} \mathbf{p}_l^T \mathbf{f}(\mathbf{x})$$

(since $f_{1l} = \mathbf{p}_l^T \mathbf{f}(\mathbf{x})$), we obtain the composite nonsmooth function (11.1) and therefore the functions $F(\mathbf{x}) = \max_{1 \leq k \leq m} f_k(\mathbf{x})$, $F(\mathbf{x}) = \|f(\mathbf{x})\|_\infty$, $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_\infty$.

- (2) Let $h : \mathbb{R}^m \rightarrow \mathbb{R}$, where $h(\mathbf{z}) = z_1 + \dots + z_m$, so $\frac{\partial}{\partial z_k} h(\mathbf{z}) = 1 > 0$, $1 \leq k \leq m$. Then function (11.2) has the form

$$F(\mathbf{x}) = \sum_{k=1}^m F_k(\mathbf{x}) = \sum_{k=1}^m \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x}) \tag{11.4}$$

(the sum of maxima). If $m_k = 2$ and $F_k(\mathbf{x}) = \max\{f_k(\mathbf{x}), -f_k(\mathbf{x})\}$, we obtain the function $F(\mathbf{x}) = \|f(\mathbf{x})\|_1$. If $m_k = 2$ and $F_k(\mathbf{x}) = \max\{f_k(\mathbf{x}), 0\}$, we obtain the function $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_1$. It follows that the expression of functions $F(\mathbf{x}) = \|f(\mathbf{x})\|_1$ and $F(\mathbf{x}) = \|f_+(\mathbf{x})\|_1$ by (11.2) contains only m summands and each summand is a maximum of two function values. Thus, this approach is much more economic than the use of formulas stated in Remark 11.1(d)–(e).

Remark 11.5 Since the functions $F_k(\mathbf{x})$, $1 \leq k \leq m$, are regular [25, Theorem 3.2.13], the function $h(\mathbf{z})$ is continuously differentiable, and $h_k = \frac{\partial}{\partial z_k} h(\mathbf{z}) > 0$, one can write [25, Theorem 3.2.9]

$$\begin{aligned} \partial F(\mathbf{x}) &= \text{conv} \sum_{k=1}^m h_k \partial F_k(\mathbf{x}) = \sum_{k=1}^m h_k \partial F_k(\mathbf{x}) \\ &= \sum_{k=1}^m h_k \text{conv}\{\mathbf{g}_{kl}(\mathbf{x}) : l \in \bar{\mathcal{L}}_k(\mathbf{x})\}, \end{aligned}$$

where $\bar{I}_k(\mathbf{x}) = \{l : 1 \leq l \leq m_k, f_{kl}(\mathbf{x}) = F_k(\mathbf{x})\}$. Thus, one has

$$\partial F(\mathbf{x}) = \sum_{k=1}^m h_k \sum_{l=1}^{m_k} \lambda_{kl} \mathbf{g}_{kl}(\mathbf{x}),$$

where for $1 \leq k \leq m$ it holds $\lambda_{kl} \geq 0$, $\lambda_{kl}(F_k(\mathbf{x}) - f_{kl}(\mathbf{x})) = 0$, $1 \leq l \leq m_k$, and $\sum_{l=1}^{m_k} \lambda_{kl} = 1$. Setting $u_{kl} = h_k \lambda_{kl}$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, we can write

$$\partial F(\mathbf{x}) = \sum_{k=1}^m \sum_{l=1}^{m_k} u_{kl} \mathbf{g}_{kl}(\mathbf{x}),$$

where for $1 \leq k \leq m$ it holds $u_{kl} \geq 0$, $u_{kl}(F_k(\mathbf{x}) - f_{kl}(\mathbf{x})) = 0$, $1 \leq l \leq m_k$, and $\sum_{l=1}^{m_k} u_{kl} = h_k$. If a point $\mathbf{x} \in \mathbb{R}^n$ is a minimum of a function $F(\mathbf{x})$, then $\mathbf{0} \in \partial F(\mathbf{x})$, so there exist multipliers u_{kl} , $1 \leq k \leq m$, $1 \leq l \leq m_k$, such that

$$\sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl} = \mathbf{0}, \quad \sum_{l=1}^{m_k} u_{kl} = h_k, \quad h_k = \frac{\partial}{\partial z_k} h(\mathbf{z}), \quad (11.5)$$

$$u_{kl} \geq 0, \quad F_k(\mathbf{x}) - f_{kl}(\mathbf{x}) \geq 0, \quad u_{kl}(F_k(\mathbf{x}) - f_{kl}(\mathbf{x})) = 0. \quad (11.6)$$

Remark 11.6 Unconstrained minimization of function (11.2) is equivalent to the nonlinear programming problem

$$\begin{cases} \text{minimize} & \tilde{F}(\mathbf{x}, \mathbf{z}) = h(\mathbf{z}) \\ \text{subject to} & f_{kl}(\mathbf{x}) \leq z_k, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k. \end{cases} \quad (11.7)$$

The condition (11.3) is sufficient for satisfying equalities $z_k = F_k(\mathbf{x})$, $1 \leq k \leq m$, at the minimum point. Denoting $c_{kl}(\mathbf{x}, \mathbf{z}) = f_{kl}(\mathbf{x}) - z_k$ and $\mathbf{a}_{kl}(\mathbf{x}, \mathbf{z}) = \nabla c_{kl}(\mathbf{x}, \mathbf{z})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, we obtain $\mathbf{a}_{kl}(\mathbf{x}, \mathbf{z}) = [\mathbf{g}_{kl}^T(\mathbf{x}), -\mathbf{e}_k^T]^T$, where $\mathbf{g}_{kl}(\mathbf{x})$ is a gradient of $f_{kl}(\mathbf{x})$ in \mathbf{x} and \mathbf{e}_k is the k -th column of the unit matrix of order m . Thus, the necessary first-order (KKT) conditions have the form

$$\mathbf{g}(\mathbf{x}, \mathbf{u}) = \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl} = \mathbf{0}, \quad \sum_{l=1}^{m_k} u_{kl} = h_k, \quad h_k = \frac{\partial}{\partial z_k} h(\mathbf{z}), \quad (11.8)$$

$$u_{kl} \geq 0, \quad z_k - f_{kl}(\mathbf{x}) \geq 0, \quad u_{kl}(z_k - f_{kl}(\mathbf{x})) = 0, \quad (11.9)$$

where u_{kl} , $1 \leq k \leq m$, $1 \leq l \leq m_k$, are Lagrange multipliers and $z_k = F_k(\mathbf{x})$. So we obtain the same necessary conditions for an extremum as in Remark 11.5.

Remark 11.7 A classical minimax problem

$$F(\mathbf{x}) = \max_{1 \leq k \leq m} f_k(\mathbf{x}) \quad (11.10)$$

can be replaced with an equivalent nonlinear programming problem

$$\begin{cases} \text{minimize} & \tilde{F}(\mathbf{x}, z) = z \\ \text{subject to} & f_k(\mathbf{x}) \leq z, \quad 1 \leq k \leq m, \end{cases} \quad (11.11)$$

and the necessary KKT conditions have the form

$$\sum_{k=1}^m \mathbf{g}_k(\mathbf{x}) u_k = \mathbf{0}, \quad \sum_{k=1}^m u_k = 1, \quad (11.12)$$

$$u_k \geq 0, \quad z - f_k(\mathbf{x}) \geq 0, \quad u_k(z - f_k(\mathbf{x})) = 0, \quad 1 \leq k \leq m. \quad (11.13)$$

Remark 11.8 Minimization of the sum of absolute values

$$F(\mathbf{x}) = \sum_{k=1}^m |f_k(\mathbf{x})| = \sum_{k=1}^m \max\{f_k^+(\mathbf{x}), f_k^-(\mathbf{x})\}, \quad (11.14)$$

where

$$f_k^+(\mathbf{x}) = f_k(\mathbf{x}), \quad f_k^-(\mathbf{x}) = -f_k(\mathbf{x})$$

can be replaced with an equivalent nonlinear programming problem

$$\begin{cases} \text{minimize} & \tilde{F}(\mathbf{x}, z) = \sum_{k=1}^m z_k \\ \text{subject to} & -z_k \leq f_k(\mathbf{x}) \leq z_k, \end{cases} \quad (11.15)$$

(there are two constraints $c_k^-(\mathbf{x}) = z_k - f_k(\mathbf{x}) \geq 0$ and $c_k^+(\mathbf{x}) = z_k + f_k(\mathbf{x}) \geq 0$ for each index $1 \leq k \leq m$) and the necessary KKT conditions have the form

$$\sum_{k=1}^m \mathbf{g}_k(\mathbf{x})(u_k^+ - u_k^-) = \mathbf{0}, \quad u_k^+ + u_k^- = 1, \quad (11.16)$$

$$u_k^+ \geq 0, \quad z_k - f_k(\mathbf{x}) \geq 0, \quad u_k^+(z_k - f_k(\mathbf{x})) = 0, \quad (11.17)$$

$$u_k^- \geq 0, \quad z_k + f_k(\mathbf{x}) \geq 0, \quad u_k^-(z_k + f_k(\mathbf{x})) = 0, \quad (11.18)$$

where $1 \leq k \leq m$. If we set $u_k = u_k^+ - u_k^-$ and use the equality $u_k^+ + u_k^- = 1$, we obtain $u_k^+ = (1 + u_k)/2$, $u_k^- = (1 - u_k)/2$. From conditions $u_k^+ \geq 0$, $u_k^- \geq 0$ the inequalities $-1 \leq u_k \leq 1$, or $|u_k| \leq 1$, follow. The condition $u_k^+ + u_k^- = 1$ implies that the numbers u_k^+ , u_k^- cannot be simultaneously zero, so either $z_k = f_k(\mathbf{x})$ or $z_k = -f_k(\mathbf{x})$, that is $z_k = |f_k(\mathbf{x})|$. If $f_k(\mathbf{x}) \neq 0$, it cannot simultaneously hold $z_k = f_k(\mathbf{x})$ and $z_k = -f_k(\mathbf{x})$, so the numbers u_k^+ , u_k^- cannot be simultaneously nonzero. Then either $u_k = u_k^+ = 1$ and $z_k = f_k(\mathbf{x})$ or $u_k = -u_k^- = -1$ and $z_k = -f_k(\mathbf{x})$, that is $u_k = f_k(\mathbf{x})/|f_k(\mathbf{x})|$. Thus, the necessary KKT conditions have the form

$$\sum_{k=1}^m \mathbf{g}_k(\mathbf{x})u_k = \mathbf{0}, \quad z_k = |f_k(\mathbf{x})|,$$

$$|u_k| \leq 1, \quad u_k = \frac{f_k(\mathbf{x})}{|f_k(\mathbf{x})|}, \quad \text{if } |f_k(\mathbf{x})| > 0. \tag{11.19}$$

Remark 11.9 Minimization of the sum of absolute values can also be reformulated so that more slack variables are used. We obtain the problem

$$\begin{cases} \text{minimize} & \tilde{F}(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^m (z_k^+ + z_k^-) \\ \text{subject to} & f_k(\mathbf{x}) = z_k^+ - z_k^-, \quad z_k^+ \geq 0, \quad z_k^- \geq 0, \end{cases} \tag{11.20}$$

where $1 \leq k \leq m$. This problem contains m general equality constraints and $2m$ simple bounds for $2m$ slack variables.

In the subsequent considerations, we will restrict ourselves to functions of the form (11.4), the sums of maxima that include most cases important for applications. In this case, it holds

$$h(\mathbf{z}) = \sum_{k=1}^m z_k, \quad \nabla h(\mathbf{z}) = \tilde{\mathbf{e}}, \quad \nabla^2 h(\mathbf{z}) = \mathbf{0}, \tag{11.21}$$

where $\tilde{\mathbf{e}} \in \mathbb{R}^m$ is a vector with unit elements. The case when $h(\mathbf{z})$ is a general function satisfying Assumption 11.3 is studied in [23].

11.2 Primal Interior Point Methods

11.2.1 Barriers and Barrier Functions

Primal interior point methods for equality constraint minimization problems are based on adding a barrier term containing constraint functions to the minimized function. A resulting barrier function, depending on a barrier parameter $0 < \mu \leq$

$\bar{\mu} < \infty$, is successively minimized on \mathbb{R}^n (without any constraints), where $\mu \downarrow 0$. Applying this approach on the problem (11.7), we obtain a barrier function

$$B_\mu(\mathbf{x}, \mathbf{z}) = h(\mathbf{z}) + \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \varphi(z_k - f_{kl}(\mathbf{x})), \quad 0 < \mu \leq \bar{\mu}, \tag{11.22}$$

where $\varphi : (0, \infty) \rightarrow \mathbb{R}$ is a barrier which satisfies the following assumption.

Assumption 11.5 *Function $\varphi(t)$, $t \in (0, \infty)$, is twice continuously differentiable, decreasing, and strictly convex, with $\lim_{t \downarrow 0} \varphi(t) = \infty$. Function $\varphi'(t)$ is increasing and strictly concave such that $\lim_{t \uparrow \infty} \varphi'(t) = 0$. For $t \in (0, \infty)$ it holds $-t\varphi'(t) \leq 1$, $t^2\varphi''(t) \leq 1$. There exist numbers $\tau > 0$ and $\underline{c} > 0$ such that for $t < \tau$ it holds*

$$-t\varphi'(t) \geq \underline{c} \tag{11.23}$$

and

$$\varphi'(t)\varphi'''(t) - \varphi''(t)^2 > 0. \tag{11.24}$$

Remark 11.10 A logarithmic barrier function

$$\varphi(t) = \log t^{-1} = -\log t, \tag{11.25}$$

is most frequently used. It satisfies Assumption 11.5 with $\underline{c} = 1$ and $\tau = \infty$ but it is not bounded from below since $\log t \uparrow \infty$ for $t \uparrow \infty$. For that reason, barriers bounded from below are sometimes used, e.g. a function

$$\varphi(t) = \log(t^{-1} + \tau^{-1}) = -\log \frac{t\tau}{t + \tau}, \tag{11.26}$$

which is bounded from below by number $\underline{\varphi} = -\log \tau$, or a function

$$\varphi(t) = -\log t, \quad 0 < t \leq \tau, \quad \varphi(t) = at^{-2} + bt^{-1} + c, \quad t \geq \tau, \tag{11.27}$$

which is bounded from below by number $\underline{\varphi} = c = -\log \tau - 3/2$, or a function

$$\varphi(t) = -\log t, \quad 0 < t \leq \tau, \quad \varphi(t) = at^{-1} + bt^{-1/2} + c, \quad t \geq \tau, \tag{11.28}$$

which is bounded from below by number $\underline{\varphi} = c = -\log \tau - 3$. Coefficients a, b, c are chosen so that function $\varphi(t)$ as well as its first and second derivatives are continuous in $t = \tau$. All these barriers satisfy Assumption 11.5 [23] (the proof of this statement is trivial for logarithmic barrier (11.25)).

Even if bounded from below barriers (11.26)–(11.28) have more advantageous theoretical properties (Assumption 11.1 can be replaced with a weaker Assump-

tion 11.2 and the proof of Lemma 11.2 below is much simpler, see [23]), algorithms using logarithmic barrier (11.26) are usually more efficient. Therefore, we will only deal with methods using the logarithmic barrier $\varphi(t) = -\log t$ in the subsequent considerations.

11.2.2 Iterative Determination of a Minimax Vector

Suppose the function $h(\mathbf{z})$ is of form (11.21). Using the logarithmic barrier $\varphi(t) = -\log t$, function (11.22) can be written as

$$B_\mu(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^m z_k - \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k - f_{kl}(\mathbf{x})), \quad 0 < \mu \leq \bar{\mu}. \quad (11.29)$$

Further, we will denote $\mathbf{g}_{kl}(\mathbf{x})$ and $G_{kl}(\mathbf{x})$ gradients and Hessian matrices of functions $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and set

$$\begin{aligned} u_{kl}(\mathbf{x}, \mathbf{z}) &= \frac{\mu}{z_k - f_{kl}(\mathbf{x})} \geq 0, \\ v_{kl}(\mathbf{x}, \mathbf{z}) &= \frac{\mu}{(z_k - f_{kl}(\mathbf{x}))^2} = \frac{1}{\mu} u_{kl}^2(\mathbf{x}, \mathbf{z}) \geq 0, \end{aligned} \quad (11.30)$$

and

$$\mathbf{u}_k(\mathbf{x}, \mathbf{z}) = \begin{bmatrix} u_{k1}(\mathbf{x}, \mathbf{z}) \\ \vdots \\ u_{km_k}(\mathbf{x}, \mathbf{z}) \end{bmatrix}, \quad \mathbf{v}_k(\mathbf{x}, \mathbf{z}) = \begin{bmatrix} v_{k1}(\mathbf{x}, \mathbf{z}) \\ \vdots \\ v_{km_k}(\mathbf{x}, \mathbf{z}) \end{bmatrix}, \quad \tilde{\mathbf{e}}_k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}.$$

Denoting by $\mathbf{g}(\mathbf{x}, \mathbf{z})$ the gradient of the function $B_\mu(\mathbf{x}, \mathbf{z})$ and $\gamma_k(\mathbf{x}, \mathbf{z}) = \frac{\partial}{\partial z_k} B_\mu(\mathbf{x}, \mathbf{z})$, the necessary conditions for an extremum of the barrier function (11.22) can be written in the form

$$\mathbf{g}(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}, \mathbf{z}) = \mathbf{0}, \quad (11.31)$$

$$\gamma_k(\mathbf{x}, \mathbf{z}) = 1 - \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}, \mathbf{z}) = 1 - \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}, \mathbf{z}) = 0, \quad 1 \leq k \leq m, \quad (11.32)$$

where $A_k(\mathbf{x}) = [\mathbf{g}_{k1}(\mathbf{x}), \dots, \mathbf{g}_{km_k}(\mathbf{x})]$, which is a system of $n + m$ nonlinear equations for unknown vectors \mathbf{x} and \mathbf{z} . These equations can be solved by the Newton method. In this case, the second derivatives of the Lagrange function (which

are the first derivatives of expressions (11.31) and (11.32)) are computed. Denoting

$$G(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^m \sum_{l=1}^{m_k} G_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}, \mathbf{z}), \quad (11.33)$$

the Hessian matrix of the Lagrange function and setting

$$U_k(\mathbf{x}, \mathbf{z}) = \text{diag}[u_{k1}(\mathbf{x}, \mathbf{z}), \dots, u_{km_k}(\mathbf{x}, \mathbf{z})],$$

$$V_k(\mathbf{x}, \mathbf{z}) = \text{diag}[v_{k1}(\mathbf{x}, \mathbf{z}), \dots, v_{km_k}(\mathbf{x}, \mathbf{z})] = \frac{1}{\mu} U_k^2(\mathbf{x}, \mathbf{z}),$$

we can write

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} \mathbf{g}(\mathbf{x}, \mathbf{z}) &= \sum_{k=1}^m \sum_{l=1}^{m_k} G_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}, \mathbf{z}) + \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) v_{kl}(\mathbf{x}, \mathbf{z}) \mathbf{g}_{kl}^T(\mathbf{x}) \\ &= G(\mathbf{x}, \mathbf{z}) + \sum_{k=1}^m A_k(\mathbf{x}) V_k(\mathbf{x}, \mathbf{z}) A_k^T(\mathbf{x}), \end{aligned} \quad (11.34)$$

$$\frac{\partial}{\partial z_k} \mathbf{g}(\mathbf{x}, \mathbf{z}) = - \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) v_{kl}(\mathbf{x}, \mathbf{z}) = -A_k(\mathbf{x}) \mathbf{v}_k(\mathbf{x}, \mathbf{z}), \quad (11.35)$$

$$\frac{\partial}{\partial \mathbf{x}} \gamma_k(\mathbf{x}, \mathbf{z}) = - \sum_{l=1}^{m_k} v_{kl}(\mathbf{x}, \mathbf{z}) \mathbf{g}_{kl}^T(\mathbf{x}) = -\mathbf{v}_k^T(\mathbf{x}, \mathbf{z}) A_k^T(\mathbf{x}), \quad (11.36)$$

$$\frac{\partial}{\partial z_k} \gamma_k(\mathbf{x}, \mathbf{z}) = \sum_{l=1}^{m_k} v_{kl}(\mathbf{x}, \mathbf{z}) = \tilde{\mathbf{e}}_k^T \mathbf{v}_k(\mathbf{x}, \mathbf{z}). \quad (11.37)$$

Using these formulas we obtain a system of linear equations describing a step of the Newton method

$$\begin{bmatrix} W(\mathbf{x}, \mathbf{z}) & -A_1(\mathbf{x}) \mathbf{v}_1(\mathbf{x}, \mathbf{z}) & \cdots & -A_m(\mathbf{x}) \mathbf{v}_m(\mathbf{x}, \mathbf{z}) \\ -\mathbf{v}_1^T(\mathbf{x}, \mathbf{z}) A_1^T(\mathbf{x}) & \tilde{\mathbf{e}}_1^T \mathbf{v}_1(\mathbf{x}, \mathbf{z}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -\mathbf{v}_m^T(\mathbf{x}, \mathbf{z}) A_m^T(\mathbf{x}) & 0 & \cdots & \tilde{\mathbf{e}}_m^T \mathbf{v}_m(\mathbf{x}, \mathbf{z}) \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta z_1 \\ \vdots \\ \Delta z_m \end{bmatrix} \quad (11.38)$$

$$= - \begin{bmatrix} \mathbf{g}(\mathbf{x}, \mathbf{z}) \\ \gamma_1(\mathbf{x}, \mathbf{z}) \\ \vdots \\ \gamma_m(\mathbf{x}, \mathbf{z}) \end{bmatrix},$$

where

$$W(\mathbf{x}, \mathbf{z}) = G(\mathbf{x}, \mathbf{z}) + \sum_{k=1}^m A_k(\mathbf{x}) V_k(\mathbf{x}, \mathbf{z}) A_k^T(\mathbf{x}). \quad (11.39)$$

Setting

$$\begin{aligned} C(\mathbf{x}, \mathbf{z}) &= [A_1(\mathbf{x})\mathbf{v}_1(\mathbf{x}, \mathbf{z}), \dots, A_m(\mathbf{x})\mathbf{v}_m(\mathbf{x}, \mathbf{z})], \\ D(\mathbf{x}, \mathbf{z}) &= \text{diag}[\tilde{\mathbf{e}}_1^T \mathbf{v}_1(\mathbf{x}, \mathbf{z}), \dots, \tilde{\mathbf{e}}_m^T \mathbf{v}_m(\mathbf{x}, \mathbf{z})] \end{aligned}$$

and $\boldsymbol{\gamma}(\mathbf{x}, \mathbf{z}) = [\gamma_1(\mathbf{x}, \mathbf{z}), \dots, \gamma_m(\mathbf{x}, \mathbf{z})]^T$, a step of the Newton method can be written in the form

$$\begin{bmatrix} W(\mathbf{x}, \mathbf{z}) & -C(\mathbf{x}, \mathbf{z}) \\ -C^T(\mathbf{x}, \mathbf{z}) & D(\mathbf{x}, \mathbf{z}) \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{z} \end{bmatrix} = - \begin{bmatrix} \mathbf{g}(\mathbf{x}, \mathbf{z}) \\ \boldsymbol{\gamma}(\mathbf{x}, \mathbf{z}) \end{bmatrix}. \quad (11.40)$$

The diagonal matrix $D(\mathbf{x}, \mathbf{z})$ is positive definite since it has positive diagonal elements.

During iterative determination of a minimax vector we know a value of the parameter μ and vectors $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{z} \in \mathbb{R}^m$ such that $z_k > F_k(\mathbf{x})$, $1 \leq k \leq m$. Using formula (11.40) we determine direction vectors $\Delta \mathbf{x}$, $\Delta \mathbf{z}$. Then, we choose a step length α so that

$$B_\mu(\mathbf{x} + \alpha \Delta \mathbf{x}, \mathbf{z} + \alpha \Delta \mathbf{z}) < B_\mu(\mathbf{x}, \mathbf{z}) \quad (11.41)$$

and $z_k + \alpha \Delta z_k > F_k(\mathbf{x} + \alpha \Delta \mathbf{x})$, $1 \leq k \leq m$. Finally, we set $\mathbf{x}_+ = \mathbf{x} + \alpha \Delta \mathbf{x}$, $\mathbf{z}_+ = \mathbf{z} + \alpha \Delta \mathbf{z}$ and determine a new value $\mu_+ < \mu$. If the matrix of system of equations (11.40) is positive definite, inequality (11.41) is satisfied for a sufficiently small value of the step length α .

Theorem 11.1 *Let the matrix $G(\mathbf{x}, \mathbf{z})$ given by (11.33) be positive definite. Then the matrix of system of equations (11.40) is positive definite.*

Proof The matrix of system of equations (11.40) is positive definite if and only if the matrix D and its Schur complement $W - CD^{-1}C^T$ are positive definite [7, Theorem 2.5.6]. The matrix D is positive definite since it has positive diagonal elements. Further, it holds

$$W - CD^{-1}C^T = G + \sum_{k=1}^m \left(A_k V_k A_k^T - A_k V_k \tilde{\mathbf{e}}_k (\tilde{\mathbf{e}}_k^T V_k \tilde{\mathbf{e}}_k)^{-1} (A_k V_k \tilde{\mathbf{e}}_k)^T \right),$$

matrices $A_k V_k A_k^T - A_k V_k \tilde{\mathbf{e}}_k (\tilde{\mathbf{e}}_k^T V_k \tilde{\mathbf{e}}_k)^{-1} (A_k V_k \tilde{\mathbf{e}}_k)^T$, $1 \leq k \leq m$, are positive semidefinite due to the Schwarz inequality and the matrix G is positive definite by the assumption. \square

11.2.3 Direct Determination of a Minimax Vector

Now we will show how to solve system of equations (11.31)–(11.32) by direct determination of a minimax vector using two-level optimization

$$\mathbf{z}(\mathbf{x}; \mu) = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} B_\mu(\mathbf{x}, \mathbf{z}), \quad (11.42)$$

and

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \hat{B}(\mathbf{x}; \mu), \quad \hat{B}(\mathbf{x}; \mu) \triangleq B_\mu(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu)). \quad (11.43)$$

The problem (11.42) serves for determination of an optimal vector $\mathbf{z}(\mathbf{x}; \mu) \in \mathbb{R}^m$. Let $\tilde{B}_\mu(\mathbf{z}) = B_\mu(\mathbf{x}, \mathbf{z})$ for a fixed chosen vector $\mathbf{x} \in \mathbb{R}^n$. The function $\tilde{B}_\mu(\mathbf{z})$ is strictly convex (as a function of a vector \mathbf{z}), since it is a sum of convex function (11.21) and strictly convex functions $-\mu \log(z_k - f_{kl}(\mathbf{x}))$, $1 \leq k \leq m$, $1 \leq l \leq m_k$. A minimum of the function $\tilde{B}_\mu(\mathbf{z})$ is its stationary point, so it is a solution of system of equations (11.32) with Lagrange multipliers (11.30). The following theorem shows that this solution exists and is unique.

Theorem 11.2 *The function $\tilde{B}_\mu(\mathbf{z}) : (F(\mathbf{x}), \infty) \rightarrow \mathbb{R}$ has a unique stationary point which is its global minimum. This stationary point is characterized by a system of equations $\boldsymbol{\gamma}(\mathbf{x}, \mathbf{z}) = \mathbf{0}$, or*

$$1 - \tilde{\mathbf{e}}_k^T \mathbf{u}_k = 1 - \sum_{l=1}^{m_k} \frac{\mu}{z_k - f_{kl}(\mathbf{x})} = 0, \quad 1 \leq k \leq m, \quad (11.44)$$

which has a unique solution $\mathbf{z}(\mathbf{x}; \mu) \in \mathbb{Z} \subset \mathbb{R}^m$ such that

$$F_k(\mathbf{x}) < F_k(\mathbf{x}) + \mu < z_k(\mathbf{x}; \mu) < F_k(\mathbf{x}) + m_k \mu \quad (11.45)$$

for $1 \leq k \leq m$.

Proof Definition 11.1 implies $f_{kl}(\mathbf{x}) \leq F_k(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, where the equality occurs for at least one index l .

(a) If (11.44) holds, then we can write

$$1 = \sum_{l=1}^{m_k} \frac{\mu}{z_k - f_{kl}(\mathbf{x})} > \frac{\mu}{z_k - F_k(\mathbf{x})} \Leftrightarrow z_k - F_k(\mathbf{x}) > \mu,$$

$$1 = \sum_{l=1}^{m_k} \frac{\mu}{z_k - f_{kl}(\mathbf{x})} < \frac{m_k \mu}{z_k - F_k(\mathbf{x})} \Leftrightarrow z_k - F_k(\mathbf{x}) < m_k \mu,$$

which proves inequalities (11.45).

(b) Since

$$\begin{aligned} \gamma_k(\mathbf{x}, F + \mu) &= 1 - \sum_{l=1}^{m_k} \frac{\mu}{\mu + F_k(\mathbf{x}) - f_{kl}(\mathbf{x})} < 1 - \frac{\mu}{\mu} = 0, \\ \gamma_k(\mathbf{x}, F + m_k\mu) &= 1 - \sum_{l=1}^{m_k} \frac{\mu}{m_k\mu + F_k(\mathbf{x}) - f_{kl}(\mathbf{x})} > 1 - \frac{m_k\mu}{m_k\mu} = 0, \end{aligned}$$

and the function $\gamma_k(\mathbf{x}, z_k)$ is continuous and decreasing in $F_k(\mathbf{x}) + \mu < z_k(\mathbf{x}; \mu) < F_k(\mathbf{x}) + m_k$ by (11.37), the equation $\gamma_k(\mathbf{x}, z_k) = 0$ has a unique solution in this interval. Since the function $\tilde{B}_\mu(z)$ is convex this solution corresponds to its global minimum. □

System (11.44) is a system of m scalar equations with localization inequalities (11.45). These scalar equations can be efficiently solved by robust methods described e.g. in [14] and [15] (details are stated in [22]). Suppose that $\mathbf{z} = \mathbf{z}(\mathbf{x}; \mu)$ and denote

$$\hat{B}(\mathbf{x}; \mu) = \sum_{k=1}^m z_k(\mathbf{x}; \mu) - \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})). \tag{11.46}$$

To find a minimum of $B_\mu(\mathbf{x}, \mathbf{z})$ in \mathbb{R}^{n+m} , it suffices to minimize $\hat{B}(\mathbf{x}; \mu)$ in \mathbb{R}^n .

Theorem 11.3 Consider the barrier function (11.46). Then

$$\nabla \hat{B}(\mathbf{x}; \mu) = \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\tilde{\mathbf{x}}; \mu), \tag{11.47}$$

$$\begin{aligned} \nabla^2 \hat{B}(\mathbf{x}; \mu) &= W(\mathbf{x}; \mu) - C(\mathbf{x}; \mu) D^{-1}(\mathbf{x}; \mu) C^T(\mathbf{x}; \mu) \\ &= G(\mathbf{x}; \mu) + \sum_{k=1}^m A_k(\mathbf{x}) V_k(\mathbf{x}; \mu) A_k^T(\mathbf{x}) \\ &\quad - \sum_{k=1}^m \frac{A_k(\mathbf{x}) V_k(\mathbf{x}; \mu) \tilde{\mathbf{e}}_k \tilde{\mathbf{e}}_k^T V_k(\mathbf{x}; \mu) A_k^T(\mathbf{x})}{\tilde{\mathbf{e}}_k^T V_k(\mathbf{x}; \mu) \tilde{\mathbf{e}}_k}, \end{aligned} \tag{11.48}$$

where $G(\mathbf{x}; \mu) = G(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu))$ and $W(\mathbf{x}; \mu)$, $C(\mathbf{x}; \mu)$, $D(\mathbf{x}; \mu)$, $U_k(\mathbf{x}; \mu)$, $V_k(\mathbf{x}; \mu) = U_k^2(\mathbf{x}; \mu)/\mu$, $1 \leq k \leq m$, are obtained by the same substitution. A solution of equation

$$\nabla^2 \hat{B}(\mathbf{x}; \mu) \Delta \mathbf{x} = -\nabla \hat{B}(\mathbf{x}; \mu) \tag{11.49}$$

is identical with $\Delta \mathbf{x}$ given by (11.40), where $\mathbf{z} = \mathbf{z}(\mathbf{x}; \mu)$ (so $\boldsymbol{\gamma}(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu)) = \mathbf{0}$).

Proof Differentiating the barrier function (11.46) and using (11.32) we obtain

$$\begin{aligned} \nabla \hat{B}(\mathbf{x}; \mu) &= \sum_{k=1}^m \frac{\partial}{\partial \mathbf{x}} z_k(\mathbf{x}; \mu) - \sum_{k=1}^m \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \left(\frac{\partial}{\partial \mathbf{x}} z_k(\mathbf{x}; \mu) - \frac{\partial}{\partial \mathbf{x}} f_{kl}(\mathbf{x}) \right) \\ &= \sum_{k=1}^m \frac{\partial}{\partial \mathbf{x}} z_k(\mathbf{x}; \mu) \left(1 - \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \right) + \sum_{k=1}^m \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \frac{\partial}{\partial \mathbf{x}} f_{kl}(\mathbf{x}) \\ &= \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) = \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu), \end{aligned}$$

where

$$u_{kl}(\mathbf{x}; \mu) = \frac{\mu}{z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})}, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k. \quad (11.50)$$

Formula (11.48) can be obtained by additional differentiation of relations (11.32) and (11.47) using (11.50). A simpler way is based on using (11.40). Since (11.32) implies $\boldsymbol{\gamma}(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu)) = \mathbf{0}$, we can substitute $\boldsymbol{\gamma} = \mathbf{0}$ into (11.40), which yields the equation

$$\left(W(\mathbf{x}, \mathbf{z}) - C(\mathbf{x}, \mathbf{z}) D^{-1}(\mathbf{x}, \mathbf{z}) C^T(\mathbf{x}, \mathbf{z}) \right) \Delta \mathbf{x} = -\mathbf{g}(\mathbf{x}, \mathbf{z}),$$

where $\mathbf{z} = \mathbf{z}(\mathbf{x}; \mu)$, that confirms validity of formulas (11.48) and (11.49) (details can be found in [22]). \square

Remark 11.11 To determine an inverse of the Hessian matrix, one can use a Woodbury formula [7, Theorem 12.1.4] which gives

$$\begin{aligned} (\nabla^2 \hat{B}(\mathbf{x}; \mu))^{-1} &= W^{-1}(\mathbf{x}; \mu) - W^{-1}(\mathbf{x}; \mu) C(\mathbf{x}; \mu) \\ &\quad \left(C^T(\mathbf{x}; \mu) W^{-1}(\mathbf{x}; \mu) C(\mathbf{x}; \mu) - D(\mathbf{x}; \mu) \right)^{-1} \\ &\quad C^T(\mathbf{x}; \mu) W^{-1}(\mathbf{x}; \mu). \end{aligned} \quad (11.51)$$

If the matrix $\nabla^2 \hat{B}(\mathbf{x}; \mu)$ is not positive definite, it can be replaced by a matrix $LL^T = \nabla^2 \hat{B}(\mathbf{x}; \mu) + E$, obtained by the Gill–Murray decomposition [10]. Note that it is more advantageous to use system of linear equations (11.40) instead of (11.49) for determination of a direction vector $\Delta \mathbf{x}$ because the system of nonlinear equations (11.44) is solved with prescribed finite precision, and thus a vector $\boldsymbol{\gamma}(\mathbf{x}, \mathbf{z})$, defined by (11.32), need not be zero.

From

$$V_k(\mathbf{x}; \mu) = \frac{1}{\mu} U_k^2(\mathbf{x}; \mu), \quad \mathbf{u}_k(\mathbf{x}; \mu) \geq 0, \quad \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}; \mu) = 1, \quad 1 \leq k \leq m,$$

it follows that $\|V_k(\mathbf{x}; \mu)\| \uparrow \infty$ if $\mu \downarrow 0$, so Hessian matrix (11.48) may be ill-conditioned if the value μ is very small. From this reason, we use a lower bound $\underline{\mu} > 0$ for μ .

Theorem 11.4 *Let Assumption 11.4 be satisfied and $\mu \geq \underline{\mu} > 0$. If $G(\mathbf{x}; \mu)$ is uniformly positive definite (if a constant \underline{G} exists such that $\mathbf{v}^T \nabla^2 G(\mathbf{x}; \mu) \mathbf{v} \geq \underline{G} \|\mathbf{v}\|^2$), then there is a number $\bar{\kappa} \geq 1$ such that $\kappa(\nabla^2 \hat{B}(\mathbf{x}; \mu)) \leq \bar{\kappa}$.*

Proof

(a) Using (11.30), (11.48), and Assumption 11.4, we obtain

$$\begin{aligned} \|\nabla^2 \hat{B}(\mathbf{x}; \mu)\| &\leq \left\| G(\mathbf{x}; \mu) + \sum_{k=1}^m A_k(\mathbf{x}) V_k(\mathbf{x}; \mu) A_k^T(\mathbf{x}) \right\| \\ &\leq \sum_{k=1}^m \sum_{l=1}^{m_k} \left(|G_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}, \mu)| + \frac{1}{\mu} \left| u_{kl}^2(\mathbf{x}; \mu) \mathbf{g}_{kl}(\mathbf{x}) \mathbf{g}_{kl}^T(\mathbf{x}) \right| \right) \\ &\leq \frac{\bar{m}}{\mu} \left(\bar{\mu} \bar{G} + \bar{g}^2 \right) \triangleq \frac{\bar{c}}{\mu} \leq \frac{\bar{c}}{\underline{\mu}}, \end{aligned} \tag{11.52}$$

because $0 \leq u_{kl}(\mathbf{x}; \mu) \leq \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}; \mu) = 1, 1 \leq k \leq m, 1 \leq l \leq m_k$, by (11.44).

(b) From the proof of Theorem 11.1 it follows that the matrix $\nabla^2 \hat{B}(\mathbf{x}; \mu) - G(\mathbf{x}; \mu)$ is positive semidefinite. Therefore,

$$\underline{\lambda}(\nabla^2 \hat{B}(\mathbf{x}; \mu)) \geq \underline{\lambda}(G(\mathbf{x}; \mu)) \geq \underline{G}.$$

(c) Since (a) implies $\bar{\lambda}(\nabla^2 \hat{B}(\mathbf{x}; \mu)) = \|\nabla^2 \hat{B}(\mathbf{x}; \mu)\| \leq \bar{c}/\underline{\mu}$, using (b) we can write

$$\kappa(\nabla^2 \hat{B}(\mathbf{x}; \mu)) = \frac{\bar{\lambda}(\nabla^2 \hat{B}(\mathbf{x}; \mu))}{\underline{\lambda}(\nabla^2 \hat{B}(\mathbf{x}; \mu))} \leq \frac{\bar{c}}{\underline{\mu} \underline{G}} \triangleq \bar{\kappa}. \tag{11.53}$$

□

Remark 11.12 If there exists a number $\bar{\kappa} > 0$ such that $\kappa(\nabla^2 \hat{B}(\mathbf{x}_i; \mu_i)) \leq \bar{\kappa}, i \in \mathbb{N}$, the direction vector $\Delta \mathbf{x}_i$, given by solving a system of equations $\nabla^2 \hat{B}(\mathbf{x}_i; \mu_i) \Delta \mathbf{x}_i = -\nabla \hat{B}(\mathbf{x}_i; \mu_i)$, satisfies the condition

$$(\Delta \mathbf{x}_i)^T \mathbf{g}(\mathbf{x}_i; \mu_i) \leq -\varepsilon_0 \|\Delta \mathbf{x}_i\| \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|, \quad i \in \mathbb{N}, \tag{11.54}$$

where $\varepsilon_0 = 1/\sqrt{\bar{\kappa}}$ and $\mathbf{g}(\mathbf{x}; \mu) = \nabla \hat{B}(\mathbf{x}; \mu)$. Then, for arbitrary numbers $0 < \varepsilon_1 \leq \varepsilon_2 < 1$ one can find a step length parameter $\alpha_i > 0$ such that for $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \Delta \mathbf{x}_i$

it holds

$$\varepsilon_1 \leq \frac{\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \hat{B}(\mathbf{x}_i; \mu_i)}{\alpha_i (\Delta \mathbf{x}_i)^T \mathbf{g}(\mathbf{x}_i; \mu_i)} \leq \varepsilon_2, \quad (11.55)$$

so there exists a number $c > 0$ such that (see [26, Section 3.2])

$$\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \hat{B}(\mathbf{x}_i; \mu_i) \leq -c \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2, \quad i \in \mathbb{N}. \quad (11.56)$$

If Assumption 11.4 is not satisfied, then only $(\Delta \mathbf{x}_i)^T \mathbf{g}(\mathbf{x}_i; \mu_i) < 0$ holds (because the matrix $\nabla^2 \hat{B}(\mathbf{x}; \mu)$ is positive definite by Theorem 11.1) and

$$\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \hat{B}(\mathbf{x}_i; \mu_i) \leq 0, \quad i \in \mathbb{N}. \quad (11.57)$$

11.2.4 Implementation

Remark 11.13 In (11.39), it is assumed that $G(\mathbf{x}, \mathbf{z})$ is the Hessian matrix of the Lagrange function. Direct computation of the matrix $G(\mathbf{x}; \mu) = G(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu))$ is usually difficult (one can use automatic differentiation as described in [13]). Thus, various approximations $G \approx G(\mathbf{x}; \mu)$ are mostly used.

- The matrix $G \approx G(\mathbf{x}; \mu)$ can be determined using differences

$$G \mathbf{w}_j = \frac{1}{\delta} \left(\sum_{k=1}^m A_k(\mathbf{x} + \delta \mathbf{w}_j) \mathbf{u}_k(\mathbf{x}; \mu) - \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu) \right).$$

The vectors \mathbf{w}_j , $1 \leq j \leq \bar{k}$, are chosen so that the number of them is as small as possible [4, 27].

- The matrix $G \approx G(\mathbf{x}; \mu)$ can be determined using the variable metric methods [17]. The vectors

$$\mathbf{d} = \mathbf{x}_+ - \mathbf{x}, \quad \mathbf{y} = \sum_{k=1}^m A_k(\mathbf{x}_+) \mathbf{u}_k(\mathbf{x}_+; \mu) - \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}_+; \mu)$$

are used for an update of G .

- If the problem is separable (i.e. $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, are functions of a small number $n_{kl} = O(1)$ of variables), one can set as in [12]

$$G = \sum_{k=1}^m \sum_{l=1}^{m_k} Z_{kl} \hat{G}_{kl} Z_{kl}^T \mathbf{u}_{kl}(\mathbf{x}, \mathbf{z}),$$

where the reduced Hessian matrices \hat{G}_{kl} are updated using the reduced vectors $\hat{\mathbf{d}}_{kl} = Z_{kl}^T(\mathbf{x}_+ - \mathbf{x})$ and $\hat{\mathbf{y}}_{kl} = Z_{kl}(\mathbf{g}_{kl}(\mathbf{x}_+) - \mathbf{g}_{kl}(\mathbf{x}))$.

Remark 11.14 The matrix $G \approx G(\mathbf{x}; \mu)$ obtained by the approach stated in Remark 11.13 can be ill-conditioned so condition (11.54) (with a chosen value $\varepsilon_0 > 0$) may not be satisfied. In this case it is possible to restart the iteration process and set $G = I$. Then $\overline{G} = 1$ and $\underline{G} = 1$ in (11.52) and (11.53), so it is a higher probability of fulfilment of condition (11.54). If the choice $G = I$ does not satisfy (11.54), we set $\Delta \mathbf{x} = -\mathbf{g}(\mathbf{x}; \mu)$ (a steepest descent direction).

An update of μ is an important part of interior point methods. Above all, $\mu \downarrow 0$ must hold, which is a main property of interior point methods. Moreover, rounding errors may cause that $z_k(\mathbf{x}; \mu) = F_k(\mathbf{x})$ when the value μ is small (because $F_k(\mathbf{x}) < z_k(\mathbf{x}; \mu) \leq F_k(\mathbf{x}) + m_k\mu$ and $F_k(\mathbf{x}) + m_k\mu \rightarrow F_k(\mathbf{x})$ if $\mu \downarrow 0$), which leads to a breakdown (division by $z_k(\mathbf{x}; \mu) - F_k(\mathbf{x}) = 0$) when computing $1/(z_k(\mathbf{x}; \mu) - F_k(\mathbf{x}))$. Therefore, we need to use a lower bound $\underline{\mu}$ for a barrier parameter (e.g. $\underline{\mu} = 10^{-8}$ when computing in double precision).

The efficiency of interior point methods also depends on the way of decreasing the value of a barrier parameter. The following heuristic procedures proved successful in practice, where \underline{g} is a suitable constant.

Procedure A

Phase 1. If $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \geq \underline{g}$, then $\mu_{i+1} = \mu_i$ (the value of a barrier parameter is unchanged).

Phase 2. If $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| < \underline{g}$, then

$$\mu_{i+1} = \max \left\{ \tilde{\mu}_{i+1}, \underline{\mu}, 10 \varepsilon_M |F(\mathbf{x}_{i+1})| \right\}, \tag{11.58}$$

where $F(\mathbf{x}_{i+1}) = F_1(\mathbf{x}_{i+1}) + \dots + F_m(\mathbf{x}_{i+1})$, ε_M is a machine precision, and

$$\tilde{\mu}_{i+1} = \min \left\{ \max\{\lambda\mu_i, \mu_i/(\sigma\mu_i + 1)\}, \max\{\|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2, 10^{-2k}\} \right\}. \tag{11.59}$$

The values $\underline{\mu} = 10^{-8}$, $\lambda = 0.85$, and $\sigma = 100$ are usually used.

Procedure B

Phase 1. If $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \geq \vartheta \mu_i$, then $\mu_{i+1} = \mu_i$ (the value of a barrier parameter is unchanged).

Phase 2. If $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 < \vartheta \mu_i$, then

$$\mu_{i+1} = \max\{\underline{\mu}, \|\mathbf{g}_i(\mathbf{x}_i; \mu_i)\|^2\}. \tag{11.60}$$

The values $\underline{\mu} = 10^{-8}$ and $\vartheta = 0.1$ are usually used.

The choice of g in Procedure **A** is not critical. We can set $\underline{g} = \infty$ but a lower value is sometimes more advantageous. Formula (11.59) requires several comments. The first argument of the minimum controls the decreasing speed of the value of a barrier parameter which is linear (a geometric sequence) for small i (the term $\lambda\mu_i$) and sublinear (a harmonic sequence) for large i (the term $\mu_i/(\sigma\mu_i + 1)$). Thus, the second argument ensuring that the value μ is small in a neighborhood of a desired solution is mainly important for large i . This situation may appear if the gradient norm $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\|$ is small even if \mathbf{x}_i is far from a solution. The idea of Procedure **B** proceeds from the fact that a barrier function $\hat{B}(\mathbf{x}; \mu)$ should be minimized with a sufficient precision for a given value of a parameter μ .

The considerations up to now are summarized in Algorithm 11.1 introduced in the Appendix. This algorithm supposes that the matrix $A(\mathbf{x})$ is sparse. If it is dense, the algorithm is simplified because there is no symbolic decomposition.

11.2.5 Global Convergence

Now we prove the global convergence of the method realized by Algorithm 11.1.

Lemma 11.1 *Let numbers $z_k(\mathbf{x}; \mu)$, $1 \leq k \leq m$, be solutions of Eq. (11.44). Then*

$$\frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) > 0, \quad 1 \leq k \leq m, \quad \frac{\partial}{\partial \mu} \hat{B}(\mathbf{x}; \mu) = - \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})).$$

Proof Differentiating (11.44) with respect to μ , one can write for $1 \leq k \leq m$

$$- \sum_{l=1}^{m_k} \frac{1}{z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})} + \sum_{l=1}^{m_k} \frac{\mu}{(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x}))^2} \frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) = 0,$$

which after multiplication of μ together with (11.30) and (11.44) gives

$$\frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) = \left(\sum_{l=1}^{m_k} \frac{\mu^2}{(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x}))^2} \right)^{-1} = \left(\sum_{l=1}^{m_k} u_{kl}^2(\mathbf{x}; \mu) \right)^{-1} > 0.$$

Differentiating the function

$$\hat{B}(\mathbf{x}; \mu) = \sum_{k=1}^m z_k(\mathbf{x}; \mu) - \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})) \quad (11.61)$$

and using (11.44) we obtain

$$\begin{aligned}
\frac{\partial}{\partial \mu} \hat{B}(\mathbf{x}; \mu) &= \sum_{k=1}^m \frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) - \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})) \\
&\quad - \sum_{k=1}^m \sum_{l=1}^{m_k} \frac{\mu}{z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})} \frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) \\
&= \frac{\partial}{\partial \mu} z_k(\mathbf{x}; \mu) \sum_{k=1}^m \left(1 - \sum_{l=1}^{m_k} \frac{\mu}{z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})} \right) \\
&\quad - \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})) \\
&= - \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})).
\end{aligned}$$

□

Lemma 11.2 *Let Assumption 11.1 be satisfied. Let $\{\mathbf{x}_i\}$ and $\{\mu_i\}$, $i \in \mathbb{N}$, be the sequences generated by Algorithm 11.1. Then the sequences $\{\hat{B}(\mathbf{x}_i; \mu_i)\}$, $\{z(\mathbf{x}_i; \mu_i)\}$, and $\{F(\mathbf{x}_i)\}$, $i \in \mathbb{N}$, are bounded. Moreover, there exists a constant $L \geq 0$ such that for $i \in \mathbb{N}$ it holds*

$$\hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) \leq \hat{B}(\mathbf{x}_{i+1}; \mu_i) + L(\mu_i - \mu_{i+1}). \quad (11.62)$$

Proof

(a) We first prove boundedness from below. Using (11.61) and Assumption 11.1, one can write

$$\begin{aligned}
\hat{B}(\mathbf{x}; \mu) - \underline{F} &= \sum_{k=1}^m z_k(\mathbf{x}; \mu) - \underline{F} - \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}; \mu) - f_{kl}(\mathbf{x})) \\
&\geq \sum_{k=1}^m (z_k(\mathbf{x}; \mu) - \underline{F} - m_k \mu \log(z_k(\mathbf{x}; \mu) - \underline{F})).
\end{aligned}$$

A convex function $\psi(t) = t - m\mu \log(t)$ has a unique minimum at the point $t = m\mu$ because $\psi'(m\mu) = 1 - m\mu/m\mu = 0$. Thus, it holds

$$\hat{B}(\mathbf{x}; \mu) \geq \underline{F} + \sum_{k=1}^m (m_k \mu - m_k \mu \log(m_k \mu))$$

$$\begin{aligned} &\geq \underline{F} + \sum_{k=1}^m \min\{0, m_k \bar{\mu}(1 - \log(m_k \bar{\mu}))\} \\ &\geq \underline{F} + \sum_{k=1}^m \min\{0, m_k \bar{\mu}(1 - \log(2m_k \bar{\mu}))\} \stackrel{\Delta}{=} \underline{B}. \end{aligned}$$

Boundedness from below of sequences $\{z(\mathbf{x}_i; \mu_i)\}$ and $\{F(\mathbf{x}_i)\}$, $i \in \mathbb{N}$, follows from inequalities (11.45) and Assumption 11.1.

(b) Now we prove boundedness from above. Similarly as in (a) we can write

$$\begin{aligned} \hat{B}(\mathbf{x}; \mu) - \underline{F} &\geq \sum_{k=1}^m \frac{z_k(\mathbf{x}; \mu) - \underline{F}}{2} \\ &\quad + \sum_{k=1}^m \left(\frac{z_k(\mathbf{x}; \mu) - \underline{F}}{2} - m_k \mu \log(z_k(\mathbf{x}; \mu) - \underline{F}) \right). \end{aligned}$$

A convex function $t/2 - m\mu \log(t)$ has a unique minimum at the point $t = 2m\mu$. Thus, it holds

$$\begin{aligned} \hat{B}(\mathbf{x}; \mu) &\geq \sum_{k=1}^m \frac{z_k(\mathbf{x}; \mu) - \underline{F}}{2} + \underline{F} + \sum_{k=1}^m \min\{0, m_k \bar{\mu}(1 - \log(2m_k \bar{\mu}))\} \\ &= \sum_{k=1}^m \frac{z_k(\mathbf{x}; \mu) - \underline{F}}{2} + \underline{B} \end{aligned}$$

or

$$\sum_{k=1}^m (z_k(\mathbf{x}; \mu) - \underline{F}) \leq 2(\hat{B}(\mathbf{x}; \mu) - \underline{B}). \tag{11.63}$$

Using the mean value theorem and Lemma 11.1, we obtain

$$\begin{aligned} &\hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \hat{B}(\mathbf{x}_{i+1}; \mu_i) \\ &= \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}_{i+1}; \tilde{\mu}_i) - f_{kl}(\mathbf{x}_{i+1}))(\mu_i - \mu_{i+1}) \\ &\leq \sum_{k=1}^m \sum_{l=1}^{m_k} \log(z_k(\mathbf{x}_{i+1}; \mu_i) - f_{kl}(\mathbf{x}_{i+1}))(\mu_i - \mu_{i+1}) \\ &\leq \sum_{k=1}^m m_k \log(z_k(\mathbf{x}_{i+1}; \mu_i) - \underline{F})(\mu_i - \mu_{i+1}), \end{aligned} \tag{11.64}$$

where $0 < \mu_{i+1} \leq \tilde{\mu}_i \leq \mu_i$. Since $\log(t) \leq t/e$ (where $e = \exp(1)$) for $t > 0$, we can write using inequalities (11.63), (11.64), and (11.45)

$$\begin{aligned} \hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \underline{B} &\leq \hat{B}(\mathbf{x}_{i+1}; \mu_i) - \underline{B} \\ &\quad + \sum_{k=1}^m m_k \log(z_k(\mathbf{x}_{i+1}; \mu_i) - \underline{F})(\mu_i - \mu_{i+1}) \\ &\leq \hat{B}(\mathbf{x}_{i+1}; \mu_i) - \underline{B} \\ &\quad + e^{-1} \sum_{k=1}^m m_k (z_k(\mathbf{x}_{i+1}; \mu_i) - \underline{F})(\mu_i - \mu_{i+1}) \\ &\leq \hat{B}(\mathbf{x}_{i+1}; \mu_i) - \underline{B} \\ &\quad + 2e^{-1} \bar{m} (\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \underline{B})(\mu_i - \mu_{i+1}) \\ &= (1 + \lambda \delta_i) (\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \underline{B}) \\ &\leq (1 + \lambda \delta_i) (\hat{B}(\mathbf{x}_i; \mu_i) - \underline{B}), \end{aligned}$$

where $\lambda = 2\bar{m}/e$ and $\delta_i = \mu_i - \mu_{i+1}$. Therefore,

$$\begin{aligned} \hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \underline{B} &\leq \prod_{j=1}^i (1 + \lambda \delta_j) (\hat{B}(\mathbf{x}_1; \mu_1) - \underline{B}) \\ &\leq \prod_{i=1}^{\infty} (1 + \lambda \delta_i) (\hat{B}(\mathbf{x}_1; \mu_1) - \underline{B}) \end{aligned} \tag{11.65}$$

and since

$$\sum_{i=1}^{\infty} \lambda \delta_i = \lambda \sum_{i=1}^{\infty} (\mu_i - \mu_{i+1}) = \lambda (\bar{\mu} - \lim_{i \uparrow \infty} \mu_i) \leq \lambda \bar{\mu}$$

the expression on the right-hand side of (11.65) is finite. Thus, the sequence $\{\hat{B}(\mathbf{x}_i; \mu_i)\}$, $i \in \mathbb{N}$, is bounded from above and the sequences $\{z(\mathbf{x}_i; \mu_i)\}$ and $\{F(\mathbf{x}_i)\}$, $i \in \mathbb{N}$, are bounded from above as well by (11.63) and (11.45).

(c) Finally, we prove formula (11.62). Using (11.64) and (11.45) we obtain

$$\begin{aligned} \hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \hat{B}(\mathbf{x}_{i+1}; \mu_i) \\ \leq \sum_{k=1}^m m_k \log(z_k(\mathbf{x}_{i+1}; \mu_i) - \underline{F})(\mu_i - \mu_{i+1}) \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{k=1}^m m_k \log(F_k(\mathbf{x}_{i+1}) + m_k \mu_i - \underline{F})(\mu_i - \mu_{i+1}) \\
&\leq \sum_{k=1}^m m_k \log(\overline{F} + m_k \overline{\mu} - \underline{F})(\mu_i - \mu_{i+1}) \\
&\stackrel{\Delta}{=} L(\mu_i - \mu_{i+1})
\end{aligned}$$

(the existence of a constant \overline{F} follows from boundedness of a sequence $\{F(\mathbf{x}_i)\}$, $i \in \mathbb{N}$), which together with (11.57) gives $\hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) \leq \hat{B}(\mathbf{x}_i; \mu_i) + L(\mu_i - \mu_{i+1})$, $i \in \mathbb{N}$. Thus, it holds

$$\hat{B}(\mathbf{x}_i; \mu_i) \leq \hat{B}(\mathbf{x}_1; \mu_1) + L(\mu_1 - \mu_i) \leq \hat{B}(\mathbf{x}_1; \mu_1) + L\overline{\mu} \stackrel{\Delta}{=} \overline{B}, \quad i \in \mathbb{N}. \quad (11.66)$$

□

The upper bounds \overline{g} and \overline{G} are not used in Lemma 11.2, so Assumption 11.4 may not be satisfied. Thus, there exists an upper bound \overline{F} (independent of \overline{g} and \overline{G}) such that $F(\mathbf{x}_i) \leq \overline{F}$ for all $i \in \mathbb{N}$. This upper bound can be used in definition of a set $\mathcal{D}_F(\overline{F})$ in Assumption 11.4.

Lemma 11.3 *Let Assumption 11.4 and the assumptions of Lemma 11.2 be satisfied. Then, if we use Procedure A or Procedure B for an update of parameter μ , the values $\{\mu_i\}$, $i \in \mathbb{N}$, form a non-decreasing sequence such that $\mu_i \downarrow 0$.*

Proof The value of parameter μ is unchanged in the first phase of Procedure A or Procedure B. Since a function $\hat{B}(\mathbf{x}; \mu)$ is continuous, bounded from below by Lemma 11.2, and since inequality (11.56) is satisfied (with $\mu_i = \mu$), it holds $\|\mathbf{g}(\mathbf{x}_i; \mu)\| \downarrow 0$ if phase 1 contains an infinite number of subsequent iterative steps [26, Section 3.2]. Thus, there exists a step (with index i) belonging to the first phase such that either $\|\mathbf{g}(\mathbf{x}_i; \mu)\| < \underline{g}$ in Procedure A or $\|\mathbf{g}(\mathbf{x}_i; \mu)\|^2 < \vartheta\mu$ in Procedure B. However, this is in contradiction with the definition of the first phase. Thus, there exists an infinite number of steps belonging to the second phase, where the value of parameter μ is decreased so that $\mu_i \downarrow 0$. □

Theorem 11.5 *Let assumptions of Lemma 11.3 be satisfied. Consider a sequence $\{\mathbf{x}_i\}$, $i \in \mathbb{N}$, generated by Algorithm 11.1, where $\underline{\delta} = \underline{\varepsilon} = \underline{\mu} = 0$. Then*

$$\begin{aligned}
\lim_{i \uparrow \infty} \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}_i) u_{kl}(\mathbf{x}_i; \mu_i) &= \mathbf{0}, \quad \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}_i; \mu_i) = 1, \\
z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i) &\geq 0, \quad u_{kl}(\mathbf{x}_i; \mu_i) \geq 0, \\
\lim_{i \uparrow \infty} u_{kl}(\mathbf{x}_i; \mu_i) (z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)) &= 0
\end{aligned}$$

for $1 \leq k \leq m$ and $1 \leq l \leq m_k$.

Proof

- (a) Equalities $\tilde{e}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i) = 1, 1 \leq k \leq m$, are satisfied by (11.44) because $\underline{\delta} = 0$. Inequalities $z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i) \geq 0$ and $u_{kl}(\mathbf{x}_i; \mu_i) \geq 0$ follow from formulas (11.45) and statement (11.50).
- (b) Relations (11.56) and (11.62) yield

$$\begin{aligned} \hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \hat{B}(\mathbf{x}_i; \mu_i) &= (\hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) - \hat{B}(\mathbf{x}_{i+1}; \mu_i)) \\ &\quad + (\hat{B}(\mathbf{x}_{i+1}; \mu_i) - \hat{B}(\mathbf{x}_i; \mu_i)) \\ &\leq L(\mu_i - \mu_{i+1}) - c \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \end{aligned}$$

and since $\lim_{i \uparrow \infty} \mu_i = 0$ (Lemma 11.3), we can write by (11.66) that

$$\begin{aligned} \underline{B} &\leq \lim_{i \uparrow \infty} \hat{B}(\mathbf{x}_{i+1}; \mu_{i+1}) \\ &\leq \hat{B}(\mathbf{x}_1; \mu_1) + L \sum_{i=1}^{\infty} (\mu_i - \mu_{i+1}) - c \sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \\ &\leq \hat{B}(\mathbf{x}_1; \mu_1) + L\bar{\mu} - c \sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 = \bar{B} - c \sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2. \end{aligned}$$

Thus, it holds

$$\sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \leq \frac{1}{c}(\bar{B} - \underline{B}) < \infty,$$

which gives $\mathbf{g}(\mathbf{x}_i; \mu_i) = \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}_i) u_{kl}(\mathbf{x}_i; \mu_i) \downarrow \mathbf{0}$.

- (c) Let indices $1 \leq k \leq m$ and $1 \leq l \leq m_k$ are chosen arbitrarily. Using (11.50) and Lemma 11.3 we obtain

$$u_{kl}(\mathbf{x}_i; \mu_i)(z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)) = \frac{\mu_i(z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i))}{z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)} = \mu_i \downarrow 0. \quad \square$$

Corollary 11.1 *Let the assumptions of Theorem 11.5 be satisfied. Then, every cluster point $\mathbf{x} \in \mathbb{R}^n$ of a sequence $\{\mathbf{x}_i\}, i \in \mathbb{N}$, satisfies necessary KKT conditions (11.8)–(11.9) where \mathbf{z} and \mathbf{u} (with elements z_k and $u_{kl}, 1 \leq k \leq m, 1 \leq l \leq m_k$) are cluster points of sequences $\{\mathbf{z}(\mathbf{x}_i; \mu_i)\}$ and $\{\mathbf{u}(\mathbf{x}_i; \mu_i)\}, i \in \mathbb{N}$.*

Now we will suppose that the values $\underline{\delta}, \underline{\varepsilon}$, and $\underline{\mu}$ are nonzero and show how a precise solution of the system of KKT equations will be after termination of computation.

Theorem 11.6 *Let the assumptions of Lemma 11.3 be satisfied. Consider a sequence $\{\mathbf{x}_i\}, i \in \mathbb{N}$, generated by Algorithm 11.1. Then, if the values $\underline{\delta} > 0$,*

$\underline{\varepsilon} > 0$, and $\underline{\mu} > 0$ are chosen arbitrarily, there exists an index $i \geq 1$ such that

$$\begin{aligned} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\| &\leq \underline{\varepsilon}, & \left| 1 - \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}_i; \mu_i) \right| &\leq \underline{\delta}, \\ z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i) &\geq 0, & u_{kl}(\mathbf{x}_i; \mu_i) &\geq 0, \\ u_{kl}(\mathbf{x}_i; \mu_i)(z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)) &\leq \underline{\mu}, \end{aligned}$$

for $1 \leq k \leq m$ and $1 \leq l \leq m_k$.

Proof Inequality $|1 - \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i)| \leq \underline{\delta}$ follows immediately from the fact that the equation $\tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i) = 1$, $1 \leq k \leq m$, is solved with precision $\underline{\delta}$. Inequalities $z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i) \geq 0$, $u_{kl}(\mathbf{x}_i; \mu_i) \geq 0$ follow from formulas (11.45) and statement (11.50) as in the proof of Theorem 11.5. Since $\mu_i \downarrow 0$ and $\mathbf{g}(\mathbf{x}_i; \mu_i) \downarrow \mathbf{0}$ by Lemma 11.3 and Theorem 11.5, there exists an index $i \geq 1$ such that $\mu_i \leq \underline{\mu}$ and $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \leq \underline{\varepsilon}$. Using (11.50) we obtain

$$u_{kl}(\mathbf{x}_i; \mu_i)(z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)) = \frac{\mu_i(z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i))}{z_k(\mathbf{x}_i; \mu_i) - f_{kl}(\mathbf{x}_i)} = \mu_i \leq \underline{\mu}.$$

□

Theorem 11.5 is a standard global convergence result. If the stopping parameters $\underline{\delta}$, $\underline{\varepsilon}$, $\underline{\mu}$ are zero, the sequence of generated points converges to the point satisfying the KKT conditions for the equivalent nonlinear programming problem. Theorem 11.6 determines a precision of the obtained solution if the stopping parameters are nonzero.

11.2.6 Special Cases

Both the simplest and most widely considered generalized minimax problem is the classical minimax problem (11.10), when $m = 1$ in (11.4) (in this case we write $m, z, \mathbf{u}, \mathbf{v}, U, V, A$ instead of $m_1, z_1, \mathbf{u}_1, \mathbf{v}_1, U_1, V_1, A_1$). For solving a classical minimax problem one can use Algorithm 11.1, where a major part of computation is very simplified. System of equations (11.38) is of order $n + 1$ and has the form

$$\begin{bmatrix} G(\mathbf{x}, z) + A(\mathbf{x})V(\mathbf{x}, z)A^T(\mathbf{x}) - A(\mathbf{x})V(\mathbf{x}, z)\tilde{\mathbf{e}} & \\ -\tilde{\mathbf{e}}^T V(\mathbf{x}, z)A^T(\mathbf{x}) & \tilde{\mathbf{e}}^T V(\mathbf{x}, z)\tilde{\mathbf{e}} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta z \end{bmatrix} = - \begin{bmatrix} \mathbf{g}(\mathbf{x}, z) \\ \gamma(\mathbf{x}, z) \end{bmatrix}, \tag{11.67}$$

where $\mathbf{g}(\mathbf{x}, z) = A(\mathbf{x})\mathbf{u}(\mathbf{x}, z)$, $\gamma(\mathbf{x}, z) = 1 - \tilde{\mathbf{e}}^T \mathbf{u}(\mathbf{x}, z)$, $V(\mathbf{x}, z) = U^2(\mathbf{x}, z)/\mu = \text{diag}[u_1^2(\mathbf{x}, z), \dots, u_m^2(\mathbf{x}, z)]/\mu$, and $u_k(\mathbf{x}, z) = \mu/(z - f_k(\mathbf{x}))$, $1 \leq k \leq m$. System

of equations (11.44) is reduced to one nonlinear equation

$$1 - \tilde{\mathbf{e}}^T \mathbf{u}(\mathbf{x}, z) = 1 - \sum_{k=1}^m \frac{\mu}{z - f_k(\mathbf{x})} = 0, \quad (11.68)$$

whose solution $z(\mathbf{x}; \mu)$ lies in the interval $F(\mathbf{x}) + \mu \leq z(\mathbf{x}; \mu) \leq F(\mathbf{x}) + m\mu$. To find this solution by robust methods from [14, 15] is not difficult. A barrier function has the form

$$\hat{B}(\mathbf{x}; \mu) = z(\mathbf{x}; \mu) - \mu \sum_{k=1}^m \log(z(\mathbf{x}; \mu) - f_k(\mathbf{x})) \quad (11.69)$$

with $\nabla \hat{B}(\mathbf{x}; \mu) = A(\mathbf{x})\mathbf{u}(\mathbf{x}; \mu)$ and

$$\nabla^2 \hat{B}(\mathbf{x}; \mu) = G(\mathbf{x}; \mu) + A(\mathbf{x})V(\mathbf{x}; \mu)A^T(\mathbf{x}) - \frac{A(\mathbf{x})V(\mathbf{x}; \mu)\tilde{\mathbf{e}}\tilde{\mathbf{e}}^T V(\mathbf{x}; \mu)A^T(\mathbf{x})}{\tilde{\mathbf{e}}^T V(\mathbf{x}; \mu)\tilde{\mathbf{e}}}.$$

If we write system (11.67) in the form

$$\begin{bmatrix} W(\mathbf{x}, z) & -\mathbf{c}(\mathbf{x}, z) \\ -\mathbf{c}^T(\mathbf{x}, z) & \delta(\mathbf{x}, z) \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta z \end{bmatrix} = - \begin{bmatrix} \mathbf{g}(\mathbf{x}, z) \\ \gamma(\mathbf{x}, z) \end{bmatrix},$$

where $W(\mathbf{x}, z) = G(\mathbf{x}, z) + A(\mathbf{x})V(\mathbf{x}, z)A^T(\mathbf{x})$, $\mathbf{c}(\mathbf{x}, z) = A(\mathbf{x})V(\mathbf{x}, z)\tilde{\mathbf{e}}$ and $\delta(\mathbf{x}, z) = \tilde{\mathbf{e}}^T V(\mathbf{x}, z)\tilde{\mathbf{e}}$, then

$$\nabla^2 \hat{B}(\mathbf{x}; \mu) = W(\mathbf{x}; \mu) - \frac{\mathbf{c}(\mathbf{x}; \mu)\mathbf{c}^T(\mathbf{x}; \mu)}{\delta(\mathbf{x}; \mu)}.$$

Since

$$\begin{bmatrix} W & -\mathbf{c} \\ -\mathbf{c}^T & \delta \end{bmatrix}^{-1} = \begin{bmatrix} W^{-1} - W^{-1}\mathbf{c}\omega^{-1}\mathbf{c}^T H^{-1} & -W^{-1}\mathbf{c}\omega^{-1} \\ -\omega^{-1}\mathbf{c}^T W^{-1} & -\omega^{-1} \end{bmatrix},$$

where $\omega = \mathbf{c}^T W^{-1}\mathbf{c} - \delta$, we can write

$$\begin{bmatrix} \Delta \mathbf{x} \\ \Delta z \end{bmatrix} = - \begin{bmatrix} W & -\mathbf{c} \\ -\mathbf{c}^T & \delta \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{g} \\ \gamma \end{bmatrix} = \begin{bmatrix} W^{-1}(\mathbf{c}\Delta z - \mathbf{g}) \\ \Delta z \end{bmatrix},$$

where

$$\Delta z = \omega^{-1}(\mathbf{c}^T W^{-1}\mathbf{g} + \gamma).$$

The matrix W is sparse if the matrix $A(\mathbf{x})$ has sparse columns. If the matrix W is not positive definite, we can use the Gill–Murray decomposition

$$W + E = LL^T, \quad (11.70)$$

where E is a positive semidefinite diagonal matrix. Then we solve equations

$$LL^T \mathbf{p} = \mathbf{g}, \quad LL^T \mathbf{q} = \mathbf{c} \quad (11.71)$$

and set

$$\Delta z = \frac{\mathbf{c}^T \mathbf{p} + \gamma}{\mathbf{c}^T \mathbf{q} - \delta}, \quad \Delta \mathbf{x} = \mathbf{q} \Delta z - \mathbf{p}. \quad (11.72)$$

If we solve the classical minimax problem, Algorithm 11.1 must be somewhat modified. In Step 2, we solve only Eq. (11.68) instead of the system of equations (11.44). In Step 4, we determine a vector $\Delta \mathbf{x}$ by solving Eq. (11.71) and using relations (11.72). In Step 4, we use the barrier function (11.69) (the nonlinear equation (11.68) must be solved at the point $\mathbf{x} + \alpha \Delta \mathbf{x}$).

Minimization of a sum of absolute values, i.e., minimization of the function (11.14) is another important generalized minimax problem. In this case, a barrier function has the form

$$\begin{aligned} B_\mu(\mathbf{x}, \mathbf{z}) &= \sum_{k=1}^m z_k - \mu \sum_{k=1}^m \log(z_k - f_k(\mathbf{x})) - \mu \sum_{k=1}^m \log(z_k + f_k(\mathbf{x})) \\ &= \sum_{k=1}^m z_k - \mu \sum_{k=1}^m \log(z_k^2 - f_k^2(\mathbf{x})), \end{aligned} \quad (11.73)$$

where $z_k > |f_k(\mathbf{x})|$, $1 \leq k \leq m$. Differentiating $B_\mu(\mathbf{x}, \mathbf{z})$ with respect to \mathbf{x} and \mathbf{z} we obtain the necessary conditions for an extremum

$$\begin{aligned} \sum_{k=1}^m \frac{2\mu f_k(\mathbf{x})}{z_k^2 - f_k^2(\mathbf{x})} \mathbf{g}_k(\mathbf{x}) &= \sum_{k=1}^m u_k(\mathbf{x}, z_k) \mathbf{g}_k(\mathbf{x}) = \mathbf{0}, \\ u_k(\mathbf{x}, z_k) &= \frac{2\mu f_k(\mathbf{x})}{z_k^2 - f_k^2(\mathbf{x})} \end{aligned} \quad (11.74)$$

and

$$1 - \frac{2\mu z_k}{z_k^2 - f_k^2(\mathbf{x})} = 1 - u_k(\mathbf{x}, z_k) \frac{z_k}{f_k(\mathbf{x})} = 0 \quad \Rightarrow \quad u_k(\mathbf{x}, z_k) = \frac{f_k(\mathbf{x})}{z_k}, \quad (11.75)$$

where $\mathbf{g}_k(\mathbf{x}) = \nabla f_k(\mathbf{x})$, $1 \leq k \leq m$, which corresponds to (11.31)–(11.32). Equations in (11.44) are quadratic of the form

$$\frac{2\mu z_k(\mathbf{x}; \mu)}{z_k^2(\mathbf{x}; \mu) - f_k^2(\mathbf{x})} = 1 \quad \Leftrightarrow \quad z_k^2(\mathbf{x}; \mu) - f_k^2(\mathbf{x}) = 2\mu z_k(\mathbf{x}; \mu), \quad (11.76)$$

where $1 \leq k \leq m$, and their solutions is given by

$$z_k(\mathbf{x}; \mu) = \mu + \sqrt{\mu^2 + f_k^2(\mathbf{x})}, \quad 1 \leq k \leq m, \quad (11.77)$$

(the second solutions of quadratic equations (11.76) do not satisfy the condition $z_k > |f_k(\mathbf{x})|$, so the obtained vector \mathbf{z} does not belong to a domain of $\tilde{B}_\mu(\mathbf{z})$). Using (11.75) and (11.77) we obtain

$$u_k(\mathbf{x}; \mu) = u_k(\mathbf{x}, z_k(\mathbf{x}; \mu)) = \frac{f_k(\mathbf{x})}{z_k(\mathbf{x}; \mu)} = \frac{f_k(\mathbf{x})}{\mu + \sqrt{\mu^2 + f_k^2(\mathbf{x})}} \quad (11.78)$$

for $1 \leq k \leq m$ and

$$\begin{aligned} \hat{B}(\mathbf{x}; \mu) &= B(\mathbf{x}, \mathbf{z}(\mathbf{x}; \mu)) = \sum_{k=1}^m z_k(\mathbf{x}; \mu) - \mu \sum_{k=1}^m \log(z_k^2(\mathbf{x}; \mu) - f_k^2(\mathbf{x})) \\ &= \sum_{k=1}^m z_k(\mathbf{x}; \mu) - \mu \sum_{k=1}^m \log(2\mu z_k(\mathbf{x}; \mu)) \\ &= \sum_{k=1}^m [z_k(\mathbf{x}; \mu) - \mu \log(z_k(\mathbf{x}; \mu))] - \mu m \log(2\mu). \end{aligned} \quad (11.79)$$

Using these expressions, we can write formulas (11.47) and (11.48) in the form

$$\nabla \hat{B}(\mathbf{x}; \mu) = \sum_{k=1}^m \mathbf{g}_k(\mathbf{x}) u_k(\mathbf{x}; \mu) \quad (11.80)$$

and

$$\nabla^2 \hat{B}(\mathbf{x}; \mu) = W(\mathbf{x}; \mu) = \sum_{k=1}^m G_k(\mathbf{x}) u_k(\mathbf{x}; \mu) + \sum_{k=1}^m \mathbf{g}_k(\mathbf{x}) v_k(\mathbf{x}; \mu) \mathbf{g}_k^T(\mathbf{x}), \quad (11.81)$$

where

$$G_k(\mathbf{x}) = \nabla^2 f_k(\mathbf{x}), \quad v_k(\mathbf{x}; \mu) = \frac{2\mu}{z_k^2(\mathbf{x}; \mu) + f_k^2(\mathbf{x})}, \quad 1 \leq k \leq m. \quad (11.82)$$

A vector $\Delta \mathbf{x} \in \mathbb{R}^n$ is determined by solving the equation

$$\nabla^2 \hat{B}(\mathbf{x}; \mu) \Delta \mathbf{x} = -\mathbf{g}(\mathbf{x}; \mu), \quad (11.83)$$

where $\mathbf{g}(\mathbf{x}; \mu) = \nabla \hat{B}(\mathbf{x}; \mu) \neq 0$. From (11.83) and (11.81) it follows

$$(\Delta \mathbf{x})^T \mathbf{g}(\mathbf{x}; \mu) = -(\Delta \mathbf{x})^T \nabla^2 \hat{B}(\mathbf{x}; \mu) \Delta \mathbf{x} \leq -(\Delta \mathbf{x})^T G(\mathbf{x}; \mu) \Delta \mathbf{x},$$

so if a matrix $G(\mathbf{x}; \mu)$ is positive definite, a matrix $\nabla \hat{B}(\mathbf{x}; \mu)$ is positive definite as well (since a diagonal matrix $V(\mathbf{x}; \mu)$ is positive definite by (11.82)) and $(\Delta \mathbf{x})^T \mathbf{g}(\mathbf{x}; \mu) < 0$ holds (a direction vector $\Delta \mathbf{x}$ is descent for a function $\hat{B}(\mathbf{x}; \mu)$).

If we minimize a sum of absolute values, Algorithm 11.1 needs to be somewhat modified. In Step 2, we solve quadratic equations (11.76) whose solutions are given by (11.77). In Step 4, we determine a vector $\Delta \mathbf{x}$ by solving Eq. (11.83), where matrix $\nabla^2 \hat{B}(\mathbf{x}; \mu)$ is given by (11.83). In Step 4, we use the barrier function (11.79).

11.3 Smoothing Methods

11.3.1 Basic Properties

Similarly as in Sect. 11.2.1 we will restrict ourselves to sums of maxima, where a mapping $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a sum of its arguments, so (11.4) holds. Smoothing methods for minimization of sums of maxima replace function (11.4) by a smoothing function

$$S(\mathbf{x}; \mu) = \sum_{k=1}^m S_k(\mathbf{x}; \mu), \quad (11.84)$$

where

$$\begin{aligned} S_k(\mathbf{x}; \mu) &= \mu \log \sum_{l=1}^{m_k} \exp \left(\frac{f_{kl}(\mathbf{x})}{\mu} \right) \\ &= F_k(\mathbf{x}) + \mu \log \sum_{l=1}^{m_k} \exp \left(\frac{f_{kl}(\mathbf{x}) - F_k(\mathbf{x})}{\mu} \right), \end{aligned} \quad (11.85)$$

depending on a smoothing parameter $0 < \mu \leq \bar{\mu}$, which is successively minimized on \mathbb{R}^n with $\mu \downarrow 0$. Since $f_{kl}(\mathbf{x}) \leq F_k(\mathbf{x})$, $1 \leq l \leq m_k$, and the equality arises for at least one index, at least one exponential function on the right-hand side of (11.85) has the value 1, so the logarithm is positive. Thus $F_k(\mathbf{x}) \leq S_k(\mathbf{x}; \mu) \leq F_k(\mathbf{x}) + \mu \log m_k$, $1 \leq k \leq m$, hold. Therefore

$$F(\mathbf{x}) \leq S(\mathbf{x}; \mu) \leq F(\mathbf{x}) + \mu \sum_{k=1}^m \log m_k, \tag{11.86}$$

so $S(\mathbf{x}; \mu) \rightarrow F(\mathbf{x})$ if $\mu \downarrow 0$.

Remark 11.15 Similarly as in Sect. 11.2.2 we will denote $\mathbf{g}_{kl}(\mathbf{x})$ and $G_{kl}(\mathbf{x})$ the gradients and Hessian matrices of functions $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and

$$\mathbf{u}_k(\mathbf{x}; \mu) = \begin{bmatrix} u_{k1}(\mathbf{x}; \mu) \\ \vdots \\ u_{km_k}(\mathbf{x}; \mu) \end{bmatrix}, \quad \tilde{\mathbf{e}}_k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix},$$

where

$$u_{kl}(\mathbf{x}; \mu) = \frac{\exp(f_{kl}(\mathbf{x})/\mu)}{\sum_{l=1}^{m_k} \exp(f_{kl}(\mathbf{x})/\mu)} = \frac{\exp((f_{kl}(\mathbf{x}) - F_k(\mathbf{x}))/\mu)}{\sum_{l=1}^{m_k} \exp((f_{kl}(\mathbf{x}) - F_k(\mathbf{x}))/\mu)}. \tag{11.87}$$

Thus, it holds $u_{kl}(\mathbf{x}; \mu) \geq 0$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and

$$\tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}; \mu) = \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) = 1. \tag{11.88}$$

Further, we denote $A_k(\mathbf{x}) = J_k^T(\mathbf{x}) = [\mathbf{g}_{k1}(\mathbf{x}), \dots, \mathbf{g}_{km_k}(\mathbf{x})]$ and $U_k(\mathbf{x}; \mu) = \text{diag}[u_{k1}(\mathbf{x}; \mu), \dots, u_{km_k}(\mathbf{x}; \mu)]$ for $1 \leq k \leq m$.

Theorem 11.7 Consider the smoothing function (11.84). Then

$$\nabla S(\mathbf{x}; \mu) = \mathbf{g}(\mathbf{x}; \mu) \tag{11.89}$$

and

$$\begin{aligned} \nabla^2 S(\mathbf{x}; \mu) &= G(\mathbf{x}; \mu) + \frac{1}{\mu} \sum_{k=1}^m A_k(\mathbf{x}) U_k(\mathbf{x}; \mu) A_k^T(\mathbf{x}) \\ &\quad - \frac{1}{\mu} \sum_{k=1}^m A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu) (A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu))^T \\ &= G(\mathbf{x}; \mu) + \frac{1}{\mu} A(\mathbf{x}) U(\mathbf{x}; \mu) A^T(\mathbf{x}) - \frac{1}{\mu} C(\mathbf{x}; \mu) C(\mathbf{x}; \mu)^T \end{aligned} \tag{11.90}$$

where $\mathbf{g}(\mathbf{x}; \mu) = \sum_{k=1}^m A_k(\mathbf{x})\mathbf{u}_k(\mathbf{x}; \mu) = A(\mathbf{x})\mathbf{u}(\mathbf{x})$ and

$$G(\mathbf{x}; \mu) = \sum_{k=1}^m G_k(\mathbf{x})\mathbf{u}_k(\mathbf{x}; \mu), \quad A(\mathbf{x}) = [A_1(\mathbf{x}), \dots, A_m(\mathbf{x})],$$

$$U(\mathbf{x}; \mu) = \text{diag}[U_1(\mathbf{x}; \mu), \dots, U_m(\mathbf{x}; \mu)],$$

$$C(\mathbf{x}; \mu) = [A_1(\mathbf{x})\mathbf{u}_1(\mathbf{x}; \mu), \dots, A_m(\mathbf{x})\mathbf{u}_m(\mathbf{x}; \mu)].$$

Proof Obviously,

$$\nabla S(\mathbf{x}; \mu) = \sum_{k=1}^m \nabla S_k(\mathbf{x}; \mu), \quad \nabla^2 S(\mathbf{x}; \mu) = \sum_{k=1}^m \nabla^2 S_k(\mathbf{x}; \mu).$$

Differentiating functions (11.85) and using (11.87) we obtain

$$\begin{aligned} \nabla S_k(\mathbf{x}; \mu) &= \frac{\mu}{\sum_{l=1}^{m_k} \exp(f_{kl}(\mathbf{x})/\mu)} \sum_{l=1}^{m_k} \frac{1}{\mu} \exp(f_{kl}(\mathbf{x})/\mu) \mathbf{g}_{kl}(\mathbf{x}) \\ &= \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) = A_k(\mathbf{x})\mathbf{u}_k(\mathbf{x}; \mu). \end{aligned} \quad (11.91)$$

Adding up these expressions yields (11.89). Further, it holds

$$\begin{aligned} \nabla u_{kl}(\mathbf{x}; \mu) &= \frac{1}{\mu} \frac{\exp(f_{kl}(\mathbf{x})/\mu)}{\sum_{l=1}^{m_k} \exp(f_{kl}(\mathbf{x})/\mu)} \mathbf{g}_{kl}(\mathbf{x}) \\ &\quad - \frac{\exp(f_{kl}(\mathbf{x})/\mu)}{(\sum_{l=1}^{m_k} \exp(f_{kl}(\mathbf{x})/\mu))^2} \sum_{l=1}^{m_k} \frac{1}{\mu} \exp(f_{kl}(\mathbf{x})/\mu) \mathbf{g}_{kl}(\mathbf{x}) \\ &= \frac{1}{\mu} u_{kl}(\mathbf{x}; \mu) \mathbf{g}_{kl}(\mathbf{x}) - \frac{1}{\mu} u_{kl}(\mathbf{x}; \mu) \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \mathbf{g}_{kl}(\mathbf{x}). \end{aligned} \quad (11.92)$$

Differentiating (11.91) and using (11.92) we obtain

$$\begin{aligned} \nabla^2 S_k(\mathbf{x}; \mu) &= \sum_{l=1}^{m_k} G_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) + \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) \nabla u_{kl}(\mathbf{x}; \mu) \\ &= G_k(\mathbf{x}; \mu) + \frac{1}{\mu} \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) \mathbf{g}_{kl}^T(\mathbf{x}) \end{aligned}$$

$$\begin{aligned}
 & - \frac{1}{\mu} \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) \left(\sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu) \right)^T \\
 & = G_k(\mathbf{x}; \mu) + \frac{1}{\mu} A_k(\mathbf{x}) U_k(\mathbf{x}; \mu) A_k^T(\mathbf{x}) \\
 & \quad - \frac{1}{\mu} A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu) (A_k(\mathbf{x}) \mathbf{u}_k(\mathbf{x}; \mu))^T,
 \end{aligned}$$

where $G_k(\mathbf{x}; \mu) = \sum_{l=1}^{m_k} G_{kl}(\mathbf{x}) u_{kl}(\mathbf{x}; \mu)$. Adding up these expressions yields (11.90). \square

Remark 11.16 Note that using (11.90) and the Schwarz inequality we obtain

$$\begin{aligned}
 \mathbf{v}^T \nabla^2 S(\mathbf{x}; \mu) \mathbf{v} & = \mathbf{v}^T G(\mathbf{x}; \mu) \mathbf{v} \\
 & \quad + \frac{1}{\mu} \sum_{k=1}^m \left(\mathbf{v}^T A_k(\mathbf{x}) U_k(\mathbf{x}; \mu) A_k^T(\mathbf{x}) \mathbf{v} - \frac{(\mathbf{v}^T A_k(\mathbf{x}) U_k(\mathbf{x}; \mu) \tilde{\mathbf{e}}_k)^2}{\tilde{\mathbf{e}}_k^T U_k(\mathbf{x}; \mu) \tilde{\mathbf{e}}_k} \right) \\
 & \geq \mathbf{v}^T G(\mathbf{x}; \mu) \mathbf{v},
 \end{aligned}$$

because $\tilde{\mathbf{e}}_k^T U_k(\mathbf{x}; \mu) \tilde{\mathbf{e}}_k = \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}; \mu) = 1$, so the Hessian matrix $\nabla^2 S(\mathbf{x}; \mu)$ is positive definite if the matrix $G(\mathbf{x}; \mu)$ is positive definite.

Using Theorem 11.7, a step of the Newton method can be written in the form $\mathbf{x}_+ = \mathbf{x} + \alpha \Delta \mathbf{x}$ where

$$\nabla^2 S(\mathbf{x}; \mu) \Delta \mathbf{x} = -\nabla S(\mathbf{x}; \mu),$$

or

$$\left(W(\mathbf{x}; \mu) - \frac{1}{\mu} C(\mathbf{x}; \mu) C^T(\mathbf{x}; \mu) \right) \Delta \mathbf{x} = -\mathbf{g}(\mathbf{x}; \mu), \tag{11.93}$$

where

$$W(\mathbf{x}; \mu) = G(\mathbf{x}; \mu) + \frac{1}{\mu} A(\mathbf{x}) U(\mathbf{x}; \mu) A^T(\mathbf{x}), \quad \mathbf{g}(\mathbf{x}; \mu) = A(\mathbf{x}) \mathbf{u}(\mathbf{x}; \mu). \tag{11.94}$$

A matrix W in (11.94) has the same structure as a matrix W in (11.48) and, by Theorem 11.7, smoothing function (11.84) has similar properties as the barrier function (11.46). Thus, one can use an algorithm that is analogous to Algorithm 11.1 and considerations stated in Remark 11.12, where $S(\mathbf{x}; \mu)$ and $\nabla^2 S(\mathbf{x}; \mu)$ are used instead of $\hat{B}(\mathbf{x}; \mu)$ and $\nabla^2 \hat{B}(\mathbf{x}; \mu)$. It means that

$$S(\mathbf{x}_{i+1}; \mu_i) - S(\mathbf{x}_i; \mu_i) \leq -c \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \quad \text{for all } i \in \mathbb{N}, \tag{11.95}$$

if Assumption 11.4 is satisfied and

$$S(\mathbf{x}_{i+1}; \mu_i) - S(\mathbf{x}_i; \mu_i) \leq 0 \quad \text{for all } i \in \mathbb{N} \quad (11.96)$$

in remaining cases.

The considerations up to now are summarized in Algorithm 11.2 introduced in the Appendix. This algorithm differs from Algorithm 11.1 in that a nonlinear equation $\tilde{\mathbf{e}}^T \mathbf{u}(\mathbf{x}; \mu) = 1$ need not be solved in Step 2 (because (11.88) follows from (11.87)), Eq. (11.93)–(11.94) instead of (11.71)–(11.72) are used in Step 4, and a barrier function $\hat{B}(\mathbf{x}; \mu)$ is replaced with a smoothing function $S(\mathbf{x}; \mu)$ in Step 6. Note that the parameter μ in (11.84) has different meaning than the same parameter in (11.46), so we could use another procedure for its update in Step 7. However, it is becoming apparent that using Procedure A or Procedure B is very efficient. On the other hand, it must be noted that using exponential functions in Algorithm 11.2 has certain disadvantages. Computation of the values of exponential functions is more time consuming than performing standard arithmetic operations and underflow may also happen (i.e. replacing nonzero values by zero values) if the value of a parameter μ is very small.

11.3.2 Global Convergence

Now we prove the global convergence of the smoothing method realized by Algorithm 11.2.

Lemma 11.4 *Choose a fixed vector $\mathbf{x} \in \mathbb{R}^n$. Then $S_k(\mathbf{x}; \mu) : (0, \infty) \rightarrow \mathbb{R}$, $1 \leq k \leq m$, are nondecreasing convex functions of $\mu > 0$ and*

$$0 \leq \log \underline{m}_k \leq \frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu) \leq \log m_k, \quad (11.97)$$

where \underline{m}_k is a number of active functions (for which $f_{kl}(\mathbf{x}) = F_k(\mathbf{x})$) and

$$\frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu) = \log \sum_{l=1}^{m_k} \exp\left(\frac{f_{kl}(\mathbf{x}) - F_k(\mathbf{x})}{\mu}\right) - \sum_{l=1}^{m_k} \left(\frac{f_{kl}(\mathbf{x}) - F_k(\mathbf{x})}{\mu}\right) u_{kl}(\mathbf{x}; \mu). \quad (11.98)$$

Proof Denoting $\varphi_{kl}(\mathbf{x}; \mu) = (f_{kl}(\mathbf{x}) - F_k(\mathbf{x}))/\mu \leq 0$, $1 \leq k \leq m$, so

$$\varphi'_{kl}(\mathbf{x}; \mu) \stackrel{\Delta}{=} \frac{\partial}{\partial \mu} \varphi_{kl}(\mathbf{x}; \mu) = -\frac{\varphi_{kl}(\mathbf{x}; \mu)}{\mu} \geq 0,$$

we can write by (11.85) that

$$S_k(\mathbf{x}; \mu) = F_k(\mathbf{x}) + \mu \log \sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu)$$

and

$$\begin{aligned} \frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu) &= \log \sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu) + \mu \frac{\sum_{l=1}^{m_k} \varphi'_{kl}(\mathbf{x}; \mu) \exp \varphi_{kl}(\mathbf{x}; \mu)}{\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu)} \\ &= \log \sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu) - \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \geq 0, \end{aligned} \quad (11.99)$$

because $\varphi_{kl}(\mathbf{x}; \mu) \leq 0$, $u_{kl}(\mathbf{x}; \mu) \geq 0$, $1 \leq k \leq m$, and $\varphi_{kl}(\mathbf{x}; \mu) = 0$ holds for at least one index. Thus, functions $S_k(\mathbf{x}; \mu)$, $1 \leq k \leq m$, are nondecreasing. Differentiating (11.87) with respect to μ we obtain

$$\begin{aligned} \frac{\partial}{\partial \mu} u_{kl}(\mathbf{x}; \mu) &= -\frac{1}{\mu} \frac{\varphi_{kl}(\mathbf{x}; \mu) \exp \varphi_{kl}(\mathbf{x}; \mu)}{\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu)} \\ &\quad + \frac{1}{\mu} \frac{\exp \varphi_{kl}(\mathbf{x}; \mu)}{\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu)} \frac{\sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) \exp \varphi_{kl}(\mathbf{x}; \mu)}{\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu)} \\ &= \frac{1}{\mu} u_{kl}(\mathbf{x}; \mu) \left(-\varphi_{kl}(\mathbf{x}; \mu) + \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \right). \end{aligned} \quad (11.100)$$

Differentiating (11.99) with respect to μ and using Eqs. (11.88) and (11.100) we can write

$$\begin{aligned} \frac{\partial^2}{\partial \mu^2} S_k(\mathbf{x}; \mu) &= -\frac{1}{\mu} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \\ &\quad + \frac{1}{\mu} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) - \frac{1}{\mu} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) \frac{\partial}{\partial \mu} u_{kl}(\mathbf{x}; \mu) \\ &= -\frac{1}{\mu} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) \frac{\partial}{\partial \mu} u_{kl}(\mathbf{x}; \mu) \\ &= \frac{1}{\mu^2} \left(\sum_{l=1}^{m_k} \varphi_{kl}^2(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \right) \left(\sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \right) \\ &\quad - \frac{1}{\mu^2} \left(\sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \right)^2 \geq 0, \end{aligned}$$

because

$$\begin{aligned} \left(\sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \right)^2 &= \left(\sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) \sqrt{u_{kl}(\mathbf{x}; \mu)} \sqrt{u_{kl}(\mathbf{x}; \mu)} \right)^2 \\ &\leq \sum_{l=1}^{m_k} \varphi_{kl}^2(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}; \mu) \end{aligned}$$

holds by the Schwarz inequality. Thus, functions $S_k(\mathbf{x}; \mu)$, $1 \leq k \leq m$, are convex, so their derivatives $\frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu)$ are nondecreasing. Obviously, it holds

$$\begin{aligned} \lim_{\mu \downarrow 0} \frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu) &= \lim_{\mu \downarrow 0} \log \sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu) - \lim_{\mu \downarrow 0} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \\ &= \log \underline{m}_k - \frac{1}{\underline{m}_k} \lim_{\mu \downarrow 0} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) \exp \varphi_{kl}(\mathbf{x}; \mu) = \log \underline{m}_k, \end{aligned}$$

because $\varphi_{kl}(\mathbf{x}; \mu) = 0$ if $f_{kl}(\mathbf{x}) = F_k(\mathbf{x})$ and $\lim_{\mu \downarrow 0} \varphi_{kl}(\mathbf{x}; \mu) = -\infty$, $\lim_{\mu \downarrow 0} \varphi_{kl}(\mathbf{x}; \mu) \exp \varphi_{kl}(\mathbf{x}; \mu) = 0$ if $f_{kl}(\mathbf{x}) < F_k(\mathbf{x})$. Similarly, it holds

$$\begin{aligned} \lim_{\mu \uparrow \infty} \frac{\partial}{\partial \mu} S_k(\mathbf{x}; \mu) &= \lim_{\mu \uparrow \infty} \log \sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}; \mu) - \lim_{\mu \uparrow \infty} \sum_{l=1}^{m_k} \varphi_{kl}(\mathbf{x}; \mu) u_{kl}(\mathbf{x}; \mu) \\ &= \log m, \end{aligned}$$

because $\lim_{\mu \uparrow \infty} \varphi_{kl}(\mathbf{x}; \mu) = 0$ and $\lim_{\mu \uparrow \infty} |u_{kl}(\mathbf{x}; \mu)| \leq 1$ for $1 \leq k \leq m$. \square

Lemma 11.5 *Let Assumptions 11.2 and 11.4 be satisfied. Then the values μ_i , $i \in \mathbb{N}$, generated by Algorithm 11.2, create a nonincreasing sequence such that $\mu_i \downarrow 0$.*

Proof Lemma 11.5 is a direct consequence of Lemma 11.3 because the same procedures for an update of a parameter μ are used and (11.95) holds. \square

Theorem 11.8 *Let the assumptions of Lemma 11.5 be satisfied. Consider a sequence $\{\mathbf{x}_i\}$ $i \in \mathbb{N}$, generated by Algorithm 11.2, where $\underline{\varepsilon} = \underline{\mu} = 0$. Then*

$$\lim_{i \uparrow \infty} \sum_{k=1}^m \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}_i; \mu_i) \mathbf{g}_{kl}(\mathbf{x}_i) = \mathbf{0}, \quad \sum_{l=1}^{m_k} u_{kl}(\mathbf{x}_i; \mu_i) = 1$$

and

$$F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i) \geq 0, \quad u_{kl}(\mathbf{x}_i; \mu_i) \geq 0, \quad \lim_{i \uparrow \infty} u_{kl}(\mathbf{x}_i; \mu_i) (F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i)) = 0$$

for $1 \leq k \leq m$ and $1 \leq l \leq m_k$.

Proof

- (a) Equations $\tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i) = 1$ for $1 \leq k \leq m$ follow from (11.88). Inequalities $F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i) \geq 0$ and $u_{kl}(\mathbf{x}_i; \mu_i) \geq 0$ for $1 \leq k \leq m$ and $1 \leq l \leq m_k$ follow from (11.4) and (11.87).
- (b) Since $S_k(\mathbf{x}; \mu)$ are nondecreasing functions of the parameter μ by Lemma 11.4 and (11.95) holds, we can write

$$\begin{aligned} \underline{F} &\leq \sum_{k=1}^m F_k(\mathbf{x}_{i+1}) \leq S(\mathbf{x}_{i+1}; \mu_{i+1}) \leq S(\mathbf{x}_{i+1}; \mu_i) \\ &\leq S(\mathbf{x}_i; \mu_i) - c \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \leq S(\mathbf{x}_1; \mu_1) - c \sum_{j=1}^i \|\mathbf{g}(\mathbf{x}_j; \mu_j)\|^2, \end{aligned}$$

where $\underline{F} = \sum_{k=1}^m \underline{F}_k$ and \underline{F}_k , $1 \leq k \leq m$, are lower bounds from Assumption 11.2. Thus, it holds

$$\underline{F} \leq \lim_{i \uparrow \infty} S(\mathbf{x}_{i+1}; \mu_{i+1}) \leq S(\mathbf{x}_1; \mu_1) - c \sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2,$$

or

$$\sum_{i=1}^{\infty} \|\mathbf{g}(\mathbf{x}_i; \mu_i)\|^2 \leq \frac{1}{c} (S(\mathbf{x}_1; \mu_1) - \underline{F}),$$

so $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \downarrow 0$, which together with inequalities $0 \leq u_{kl}(\mathbf{x}_i; \mu_i) \leq 1$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, gives $\lim_{i \uparrow \infty} u_{kl}(\mathbf{x}_i; \mu_i) \mathbf{g}_{kl}(\mathbf{x}_i) = \mathbf{0}$.

- (c) Let indices $1 \leq k \leq m$ and $1 \leq l \leq m_k$ be chosen arbitrarily. Using (11.87) we get

$$\begin{aligned} 0 &\leq u_{kl}(\mathbf{x}_i; \mu_i) (F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i)) = -\mu_i \frac{\varphi_{kl}(\mathbf{x}_i; \mu_i) \exp \varphi_{kl}(\mathbf{x}_i; \mu_i)}{\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}_i; \mu_i)} \\ &\leq -\mu_i \varphi_{kl}(\mathbf{x}_i; \mu_i) \exp \varphi_{kl}(\mathbf{x}_i; \mu_i) \leq \frac{\mu_i}{e}, \end{aligned}$$

where $\varphi_{kl}(\mathbf{x}_i; \mu_i)$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, are functions used in the proof of Lemma 11.4, because

$$\sum_{l=1}^{m_k} \exp \varphi_{kl}(\mathbf{x}_i; \mu_i) \geq 1$$

and the function $t \exp t$ attains its minimal value $-1/e$ at the point $t = -1$. Since $\mu_i \downarrow 0$, we obtain $u_{kl}(\mathbf{x}_i; \mu_i)(F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i)) \downarrow 0$.

□

Corollary 11.2 *Let the assumptions of Theorem 11.8 be satisfied. Then every cluster point $\mathbf{x} \in \mathbb{R}^n$ of a sequence $\{\mathbf{x}_i\}$, $i \in \mathbb{N}$, satisfies the necessary KKT conditions (11.5) and (11.6), where \mathbf{u} (with elements u_k , $1 \leq k \leq m$) is a cluster point of a sequence $\{\mathbf{u}(\mathbf{x}_i; \mu_i)\}$, $i \in \mathbb{N}$.*

Now we will suppose that the values $\underline{\varepsilon}$ and $\underline{\mu}$ are nonzero and show how a precise solution of the system of KKT equations will be after termination of computation of Algorithm 11.2.

Theorem 11.9 *Let the assumptions of Theorem 11.5 be satisfied and let $\{\mathbf{x}_i\}$, $i \in \mathbb{N}$, be a sequence generated by Algorithm 11.2. Then, if the values $\underline{\varepsilon} > 0$ and $\underline{\mu} > 0$ are chosen arbitrarily, there exists an index $i \geq 1$ such that*

$$\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \leq \underline{\varepsilon}, \quad \tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i) = 1, \quad 1 \leq k \leq m,$$

and

$$F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i) \geq 0, \quad u_{kl}(\mathbf{x}_i; \mu_i) \geq 0, \quad u_{kl}(\mathbf{x}_i; \mu_i)(F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i)) \leq \frac{\underline{\mu}}{e}$$

for all $1 \leq k \leq m$ and $1 \leq l \leq m_k$.

Proof Equalities $\tilde{\mathbf{e}}_k^T \mathbf{u}_k(\mathbf{x}_i; \mu_i) = 1$, $1 \leq k \leq m$, follow from (11.88). Inequalities $F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i) \geq 0$ and $u_{kl}(\mathbf{x}_i; \mu_i) \geq 0$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, follow from (11.10) and (11.87). Since $\mu_i \downarrow 0$ holds by Lemma 11.5 and $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \downarrow 0$ holds by Theorem 11.8, there exists an index $i \geq 1$ such that $\mu_i \leq \underline{\mu}$ and $\|\mathbf{g}(\mathbf{x}_i; \mu_i)\| \leq \underline{\varepsilon}$. By (11.87), as in the proof of Theorem 11.8, one can write

$$u_{kl}(\mathbf{x}_i; \mu_i)(F_k(\mathbf{x}_i) - f_{kl}(\mathbf{x}_i)) \leq -\mu_i \varphi_{kl}(\mathbf{x}_i; \mu_i) \exp \varphi_{kl}(\mathbf{x}_i; \mu_i) \leq \frac{\mu_i}{e} \leq \frac{\underline{\mu}}{e}$$

for $1 \leq k \leq m$ and $1 \leq l \leq m_k$. □

Theorems 11.8 and 11.9 have the same meaning as Theorems 11.5 and 11.6 introduced in Sect. 11.2.5.

11.3.3 Special Cases

Both the simplest and most widely considered generalized minimax problem is the classical minimax problem (11.10), when $m = 1$ in (11.4) (in this case we write m and z instead of m_1 and z_1). For solving a classical minimax problem one can use Algorithm 11.2, where a major part of computation is very simplified. A step of the

Newton method can be written in the form $\mathbf{x}_+ = \mathbf{x} + \alpha \Delta \mathbf{x}$ where

$$\nabla^2 S(\mathbf{x}; \mu) \Delta \mathbf{x} = -\nabla S(\mathbf{x}; \mu),$$

or

$$\left(W(\mathbf{x}; \mu) - \frac{1}{\mu} \mathbf{g}(\mathbf{x}; \mu) \mathbf{g}^T(\mathbf{x}; \mu) \right) \Delta \mathbf{x} = -\mathbf{g}(\mathbf{x}; \mu), \quad (11.101)$$

where

$$W(\mathbf{x}; \mu) = G(\mathbf{x}; \mu) + \frac{1}{\mu} A(\mathbf{x}) U(\mathbf{x}; \mu) A^T(\mathbf{x}), \quad \mathbf{g}(\mathbf{x}; \mu) = A(\mathbf{x}) \mathbf{u}(\mathbf{x}; \mu). \quad (11.102)$$

Since

$$\left(W - \frac{1}{\mu} \mathbf{g} \mathbf{g}^T \right)^{-1} = W^{-1} + \frac{W^{-1} \mathbf{g} \mathbf{g}^T W^{-1}}{\mu - \mathbf{g}^T W^{-1} \mathbf{g}}$$

holds by the Sherman–Morrison formula, the solution of system of equations (11.101) can be written in the form

$$\Delta \mathbf{x} = \frac{\mu}{\mathbf{g}^T W^{-1} \mathbf{g} - \mu} W^{-1} \mathbf{g}. \quad (11.103)$$

If a matrix W is not positive definite, it may be replaced with a matrix $LL^T = W + E$ obtained by the Gill–Murray decomposition described in [10]. Then, we solve an equation

$$LL^T \mathbf{p} = \mathbf{g}, \quad (11.104)$$

and set

$$\Delta \mathbf{x} = \frac{\mu}{\mathbf{g}^T \mathbf{p} - \mu} \mathbf{p}. \quad (11.105)$$

Minimization of a sum of absolute values, i.e., minimization of the function (11.14) is another important generalized minimax problem. In this case, a smoothing function has the form

$$\begin{aligned} S(\mathbf{x}; \mu) &= F(\mathbf{x}) \\ &+ \mu \sum_{k=1}^m \log \left(\exp \left(-\frac{|f_k(\mathbf{x})| - f_k^+(\mathbf{x})}{\mu} \right) + \exp \left(-\frac{|f_k(\mathbf{x})| - f_k^-(\mathbf{x})}{\mu} \right) \right) \\ &= \sum_{k=1}^m |f_k(\mathbf{x})| + \mu \sum_{k=1}^m \log \left(1 + \exp \left(-\frac{2|f_k(\mathbf{x})|}{\mu} \right) \right), \end{aligned}$$

because $f_k^+(\mathbf{x}) = |f_k(\mathbf{x})|$ if $f_k(\mathbf{x}) \geq 0$ and $f_k^-(\mathbf{x}) = |f_k(\mathbf{x})|$ if $f_k(\mathbf{x}) \leq 0$, and by Theorem 11.7 we have

$$\begin{aligned}\nabla S(\mathbf{x}; \mu) &= \sum_{k=1}^m (\mathbf{g}_k^+ u_k^+ + \mathbf{g}_k^- u_k^-) = \sum_{k=1}^m \mathbf{g}_k (u_k^+ - u_k^-) = \sum_{k=1}^m \mathbf{g}_k u_k = \mathbf{g}(\mathbf{x}; \mu), \\ \nabla^2 S(\mathbf{x}; \mu) &= \sum_{k=1}^m G_k (u_k^+ - u_k^-) + \frac{1}{\mu} \sum_{k=1}^m \mathbf{g}_k \mathbf{g}_k^T (u_k^+ + u_k^-) \\ &\quad - \frac{1}{\mu} \sum_{k=1}^m \mathbf{g}_k \mathbf{g}_k^T (u_k^+ - u_k^-)^2 = G(\mathbf{x}; \mu) + \frac{1}{\mu} \sum_{k=1}^m \mathbf{g}_k \mathbf{g}_k^T (1 - u_k^2),\end{aligned}$$

(because $u_k^+ + u_k^- = 1$), where $\mathbf{g}_k = \mathbf{g}_k(\mathbf{x})$,

$$\begin{aligned}u_k &= u_k^+ - u_k^- = \frac{\exp\left(-\frac{|f_k(\mathbf{x})| - f_k^+(\mathbf{x})}{\mu}\right) - \exp\left(-\frac{|f_k(\mathbf{x})| - f_k^-(\mathbf{x})}{\mu}\right)}{\exp\left(-\frac{|f_k(\mathbf{x})| - f_k^+(\mathbf{x})}{\mu}\right) + \exp\left(-\frac{|f_k(\mathbf{x})| - f_k^-(\mathbf{x})}{\mu}\right)} \\ &= \frac{1 - \exp\left(-\frac{2|f_k(\mathbf{x})|}{\mu}\right)}{1 + \exp\left(-\frac{2|f_k(\mathbf{x})|}{\mu}\right)} \text{sign}(f_k(\mathbf{x})),\end{aligned}$$

and

$$1 - u_k^2 = \frac{4 \exp\left(-\frac{2|f_k(\mathbf{x})|}{\mu}\right)}{\left(1 + \exp\left(-\frac{2|f_k(\mathbf{x})|}{\mu}\right)\right)^2},$$

and where $\text{sign}(f_k(\mathbf{x}))$ is a sign of a function $f_k(\mathbf{x})$.

11.4 Primal-Dual Interior Point Methods

11.4.1 Basic Properties

Primal interior point methods for solving nonlinear programming problems profit from the simplicity of obtaining and keeping a point in the interior of the feasible set (for generalized minimax problems, it suffices to set $z_k > F_k(\mathbf{x})$, $1 \leq k \leq m$). Minimization of a barrier function without constraints and a direct computation of multipliers u_{kl} , $1 \leq k \leq m$, $1 \leq l \leq m_k$, are basic features of these methods. Primal-dual interior point methods are intended for solving general nonlinear programming

problems, where it is usually impossible to assure validity of constraints. These methods guarantee feasibility of points by adding slack variables, which appear in a barrier term added to the objective function. Positivity of the slack variables is assured algorithmically (by a step length selection). Minimization of a barrier function with equality constraints and an iterative computation of the Lagrange multipliers (dual variables) are the main features of primal-dual interior point methods.

Consider function (11.4). As is mentioned in the introduction, minimization of this function is equivalent to the nonlinear programming problem

$$\begin{cases} \text{minimize} & \sum_{k=1}^m z_k \\ \text{subject to} & f_{kl}(\mathbf{x}) \leq z_k, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k. \end{cases} \quad (11.106)$$

Using slack variables $s_{kl} > 0$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and a barrier function

$$B_\mu(\mathbf{x}, \mathbf{z}, \mathbf{s}) = \sum_{k=1}^m z_k - \mu \sum_{k=1}^m \sum_{l=1}^{m_k} \log(s_{kl}), \quad (11.107)$$

a solving of the problem (11.106) can be transformed to a successive solving of problems

$$\begin{cases} \text{minimize} & B_\mu(\mathbf{x}, \mathbf{z}, \mathbf{s}) \\ \text{subject to} & f_{kl}(\mathbf{x}) + s_{kl} - z_k = 0, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k, \end{cases} \quad (11.108)$$

where $\mu \downarrow 0$. Necessary conditions for an extremum of the problem (11.108) have the form

$$\begin{aligned} \mathbf{g}(\mathbf{x}, \mathbf{u}) &= \sum_{k=1}^m \sum_{l=1}^{m_k} \mathbf{g}_{kl}(\mathbf{x}) u_{kl} = \mathbf{0}, \\ 1 - \sum_{l=1}^{m_k} u_{kl} &= 0, \quad 1 \leq k \leq m, \\ u_{kl} s_{kl} - \mu &= 0, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k, \\ f_{kl}(\mathbf{x}) + s_{kl} - z_k &= 0, \quad 1 \leq k \leq m, \quad 1 \leq l \leq m_k, \end{aligned}$$

which is $n + m + 2\bar{m}$ equations for $n + m + 2\bar{m}$ unknowns (vectors \mathbf{x} , $\mathbf{z} = [z_k]$, $\mathbf{s} = [s_{kl}]$, $\mathbf{u} = [u_{kl}]$, $1 \leq k \leq m$, $1 \leq l \leq m_k$), where $\bar{m} = m_1 + \dots + m_m$. Denote

$A(\mathbf{x}) = [A_1(\mathbf{x}), \dots, A_m(\mathbf{x})]$, $\mathbf{f} = [f_{kl}]$, $S = \text{diag}[s_{kl}]$, $U = \text{diag}[u_{kl}]$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and

$$E = \begin{bmatrix} \tilde{\mathbf{e}}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{e}}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \tilde{\mathbf{e}}_m \end{bmatrix}, \quad \tilde{\mathbf{e}} = \begin{bmatrix} \tilde{e}_1 \\ \tilde{e}_2 \\ \vdots \\ \tilde{e}_m \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

(matrices $A_k(\mathbf{x})$, vectors $\tilde{\mathbf{e}}_k$, and numbers z_k , $1 \leq k \leq m$, are defined in Sect. 11.2.2). Applying the Newton method to this system of nonlinear equations, we obtain a system of linear equations for increments (direction vectors) $\Delta \mathbf{x}$, $\Delta \mathbf{z}$, $\Delta \mathbf{s}$, $\Delta \mathbf{u}$. After arrangement and elimination

$$\Delta \mathbf{s} = -U^{-1}S(\mathbf{u} + \Delta \mathbf{u}) + \mu S^{-1}\tilde{\mathbf{e}}, \tag{11.109}$$

this system has the form

$$\begin{bmatrix} G(\mathbf{x}, \mathbf{u}) & \mathbf{0} & A(\mathbf{x}) \\ \mathbf{0} & \mathbf{0} & -E^T \\ A^T(\mathbf{x}) & -E & -U^{-1}S \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{z} \\ \Delta \mathbf{u} \end{bmatrix} = - \begin{bmatrix} \mathbf{g}(\mathbf{x}, \mathbf{u}) \\ \tilde{\mathbf{e}} - E^T \mathbf{u} \\ \mathbf{f}(\mathbf{x}) - E\mathbf{z} + \mu U^{-1}\tilde{\mathbf{e}} \end{bmatrix}, \tag{11.110}$$

where $G(\mathbf{x}, \mathbf{u}) = \sum_{k=1}^m \sum_{l=1}^{m_k} G_{kl}(\mathbf{x})u_{kl}$. Vector $\tilde{\mathbf{e}}$ in the equation $\tilde{\mathbf{e}} - E^T \mathbf{u} = \mathbf{0}$ has unit elements, but its dimension is different from the dimension of a vector $\tilde{\mathbf{e}}$ in (11.109).

For solving this linear system, we cannot advantageously use the structure of a generalized minimax problem (because substituting $z_k = F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x})$ we would obtain a nonsmooth problem whose solution is much more difficult). Therefore, we need to deal with a general nonlinear programming problem. To simplify subsequent considerations, we use the notation $\tilde{\mathbf{x}} = [\mathbf{x}^T, \mathbf{z}^T]^T$,

$$\tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) = \begin{bmatrix} \mathbf{g}(\mathbf{x}, \mathbf{u}) \\ \tilde{\mathbf{e}} - E^T \mathbf{u} \end{bmatrix}, \quad \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) = \begin{bmatrix} G(\mathbf{x}, \mathbf{u}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{A}(\tilde{\mathbf{x}}) = \begin{bmatrix} A(\mathbf{x}) \\ -E^T \end{bmatrix}, \tag{11.111}$$

and write (11.110) in the form

$$\begin{bmatrix} \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) & \tilde{A}(\tilde{\mathbf{x}}) \\ \tilde{A}^T(\tilde{\mathbf{x}}) & -U^{-1}S \end{bmatrix} \begin{bmatrix} \Delta \tilde{\mathbf{x}} \\ \Delta \mathbf{u} \end{bmatrix} = - \begin{bmatrix} \tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) \\ \mathbf{c}(\tilde{\mathbf{x}}) + \mu U^{-1}\tilde{\mathbf{e}} \end{bmatrix}, \tag{11.112}$$

where $\mathbf{c}(\tilde{\mathbf{x}}) = \mathbf{f}(\mathbf{x}) - E\mathbf{z}$. This system of equations is more advantageous against systems (11.49) and (11.93) in that its matrix does not depend on the barrier parameter μ , so it is not necessary to use a lower bound $\underline{\mu}$. On the other hand, system (11.112) has a dimension $n + m + \bar{m}$, while systems (11.49) and (11.93)

have dimensions n . It would be possible to eliminate the vector $\Delta \mathbf{u}$, so the resulting system

$$(\tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) + \tilde{A}(\tilde{\mathbf{x}})M^{-1}\tilde{A}^T(\tilde{\mathbf{x}}))\Delta\tilde{\mathbf{x}} = -\tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) - \tilde{A}(\tilde{\mathbf{x}})(M^{-1}\mathbf{c}(\tilde{\mathbf{x}}) + \mu S^{-1}\tilde{\mathbf{e}}), \tag{11.113}$$

where $M = U^{-1}S$, would have dimension $n + m$ (i.e., $n + 1$ for classical minimax problems). Nevertheless, as follows from the equation $u_{kl}s_{kl} = \mu$, either $u_{kl} \downarrow 0$ or $s_{kl} \downarrow 0$ if $\mu \downarrow 0$, so some elements of a matrix M^{-1} may tend to infinity, which increases the condition number of system (11.113). Conversely, the solution of Eq. (11.112) is easier if the elements of a matrix M are small (if $M = 0$, we obtain the saddle point system, which can be solved by efficient iterative methods [1, 18]). Therefore, it is advantageous to split the constraints to active with $s_{kl} \leq \tilde{\epsilon}u_{kl}$ (we denote active quantities by $\hat{\mathbf{c}}(\tilde{\mathbf{x}})$, $\hat{A}(\tilde{\mathbf{x}})$, $\hat{\mathbf{s}}$, $\Delta\hat{\mathbf{s}}$, \hat{S} , $\hat{\mathbf{u}}$, $\Delta\hat{\mathbf{u}}$, \hat{U} , $\hat{M} = \hat{U}^{-1}\hat{S}$) and inactive with $s_{kl} > \tilde{\epsilon}u_{kl}$ (we denote inactive quantities by $\check{\mathbf{c}}(\tilde{\mathbf{x}})$, $\check{A}(\tilde{\mathbf{x}})$, $\check{\mathbf{s}}$, $\Delta\check{\mathbf{s}}$, \check{S} , $\check{\mathbf{u}}$, $\Delta\check{\mathbf{u}}$, \check{U} , $\check{M} = \check{U}^{-1}\check{S}$). Eliminating inactive equations from (11.112) we obtain

$$\Delta\check{\mathbf{u}} = \check{M}^{-1}(\check{\mathbf{c}}(\tilde{\mathbf{x}}) + \check{A}(\tilde{\mathbf{x}})^T \Delta\tilde{\mathbf{x}}) + \mu\check{S}^{-1}\tilde{\mathbf{e}} \tag{11.114}$$

and

$$\begin{bmatrix} \hat{G}(\tilde{\mathbf{x}}, \mathbf{u}) & \hat{A}(\tilde{\mathbf{x}}) \\ \hat{A}^T(\tilde{\mathbf{x}}) & -\hat{M} \end{bmatrix} \begin{bmatrix} \Delta\tilde{\mathbf{x}} \\ \Delta\hat{\mathbf{u}} \end{bmatrix} = - \begin{bmatrix} \hat{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) \\ \hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \mu\hat{U}^{-1}\tilde{\mathbf{e}} \end{bmatrix}, \tag{11.115}$$

where

$$\begin{aligned} \hat{G}(\tilde{\mathbf{x}}, \mathbf{u}) &= G(\tilde{\mathbf{x}}, \mathbf{u}) + \check{A}(\tilde{\mathbf{x}})\check{M}^{-1}\check{A}^T(\tilde{\mathbf{x}}), \\ \hat{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) &= \mathbf{g}(\tilde{\mathbf{x}}, \mathbf{u}) + \check{A}(\tilde{\mathbf{x}})(\check{M}^{-1}\check{\mathbf{c}}(\tilde{\mathbf{x}}) + \mu\check{S}^{-1}\tilde{\mathbf{e}}), \end{aligned}$$

and $\hat{M} = \hat{U}^{-1}\hat{S}$ is a diagonal matrix of order \hat{m} , where $0 \leq \hat{m} \leq \bar{m}$ is the number of active constraints. Substituting (11.114) into (11.109) we can write

$$\Delta\hat{\mathbf{s}} = -\hat{M}(\hat{\mathbf{u}} + \Delta\hat{\mathbf{u}}) + \mu\hat{U}^{-1}\tilde{\mathbf{e}}, \quad \Delta\check{\mathbf{s}} = -(\check{\mathbf{c}} + \check{A}^T \Delta\tilde{\mathbf{x}} + \check{\mathbf{s}}). \tag{11.116}$$

The matrix of the linear system (11.115) is symmetric, but indefinite, so its Choleski decomposition cannot be determined. In this case, we use either dense [3] or sparse [6] Bunch–Parlett decomposition for solving this system. System (11.115) (especially if it is large and sparse) can be efficiently solved by iterative conjugate gradient method with indefinite preconditioner [20]. If the vectors $\Delta\tilde{\mathbf{x}}$ and $\Delta\hat{\mathbf{u}}$ are solutions of system (11.115), we determine vector $\Delta\check{\mathbf{u}}$ by (11.114) and vectors $\Delta\hat{\mathbf{s}}$, $\Delta\check{\mathbf{s}}$ by (11.116).

Having vectors $\Delta\tilde{\mathbf{x}}$, $\Delta\mathbf{s}$, $\Delta\mathbf{u}$, we need to determine a step length $\alpha > 0$ and set

$$\tilde{\mathbf{x}}_+ = \tilde{\mathbf{x}} + \alpha\Delta\tilde{\mathbf{x}}, \quad \mathbf{s}_+ = \mathbf{s}(\alpha), \quad \mathbf{u}_+ = \mathbf{u}(\alpha), \tag{11.117}$$

where $s(\alpha)$ and $\mathbf{u}(\alpha)$ are vector functions such that $s(\alpha) > 0$, $s'(0) = \Delta s$ and $\mathbf{u}(\alpha) > 0$, $\mathbf{u}'(0) = \Delta \mathbf{u}$. This step is not trivial, because we need to decrease both the value of the barrier function $\tilde{B}_\mu(\tilde{\mathbf{x}}, s) = B_\mu(\mathbf{x}, \mathbf{z}, s)$ and the norm of constraints $\|\mathbf{c}(\tilde{\mathbf{x}})\|$, and also to assure positivity of vectors s and \mathbf{u} . We can do this in several different ways: using either the augmented Lagrange function [20, 21] or a bi-criterial filter [9, 29] or a special algorithm [11, 16]. In this section, we confine our attention to the augmented Lagrange function which has (for the problem (11.106)) the form

$$P(\alpha) = \tilde{B}_\mu(\tilde{\mathbf{x}} + \alpha \Delta \tilde{\mathbf{x}}, s(\alpha)) + (\mathbf{u} + \Delta \mathbf{u})^T (\mathbf{c}(\tilde{\mathbf{x}} + \alpha \Delta \tilde{\mathbf{x}}) + s(\alpha)) + \frac{\sigma}{2} \|\mathbf{c}(\tilde{\mathbf{x}} + \alpha \Delta \tilde{\mathbf{x}}) + s(\alpha)\|^2, \quad (11.118)$$

where $\sigma \geq 0$ is a penalty parameter. The following theorem, whose proof is given in [20], holds.

Theorem 11.10 *Let $s > 0$, $\mathbf{u} > 0$ and let vectors $\Delta \tilde{\mathbf{x}}$, $\Delta \hat{\mathbf{u}}$ be solutions of the linear system*

$$\begin{bmatrix} \hat{G}(\tilde{\mathbf{x}}, \mathbf{u}) & \hat{A}(\tilde{\mathbf{x}}) \\ \hat{A}^T(\tilde{\mathbf{x}}) & -\hat{M} \end{bmatrix} \begin{bmatrix} \Delta \tilde{\mathbf{x}} \\ \Delta \hat{\mathbf{u}} \end{bmatrix} + \begin{bmatrix} \hat{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) \\ \hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \mu \hat{U}^{-1} \hat{\mathbf{e}} \end{bmatrix} = \begin{bmatrix} \mathbf{r} \\ \hat{\mathbf{r}} \end{bmatrix}, \quad (11.119)$$

where \mathbf{r} and $\hat{\mathbf{r}}$ are residual vectors, and let vectors $\Delta \check{\mathbf{u}}$ and Δs be determined by (11.114) and (11.116). Then

$$P'(0) = -(\Delta \tilde{\mathbf{x}})^T \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) \Delta \tilde{\mathbf{x}} - (\Delta s)^T M^{-1} \Delta s - \sigma \|\mathbf{c}(\tilde{\mathbf{x}}) + s\|^2 + (\Delta \tilde{\mathbf{x}})^T \mathbf{r} + \sigma (\hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \hat{s})^T \hat{\mathbf{r}}. \quad (11.120)$$

If

$$\sigma > - \frac{(\Delta \tilde{\mathbf{x}})^T \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) \Delta \tilde{\mathbf{x}} + (\Delta s)^T M^{-1} \Delta s}{\|\mathbf{c}(\tilde{\mathbf{x}}) + s\|^2} \quad (11.121)$$

and if system (11.115) is solved in such a way that

$$(\Delta \tilde{\mathbf{x}})^T \mathbf{r} + \sigma (\hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \hat{s})^T \hat{\mathbf{r}} < (\Delta \tilde{\mathbf{x}})^T \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) \Delta \tilde{\mathbf{x}} + (\Delta s)^T M^{-1} \Delta s + \sigma (\|\mathbf{c}(\tilde{\mathbf{x}}) + s\|^2), \quad (11.122)$$

then $P'(0) < 0$.

Inequality (11.122) is significant only if linear system (11.115) is solved iteratively and residual vectors \mathbf{r} and $\hat{\mathbf{r}}$ are nonzero. If these vectors are zero, then (11.122) follows immediately from (11.121). Inequality (11.121) serves for determination of a penalty parameter, which should be as small as possible. If the matrix $\tilde{G}(\tilde{\mathbf{x}}, \mathbf{u})$ is positive semidefinite, then the right-hand side of (11.121) is negative and we can choose $\sigma = 0$.

11.4.2 Implementation

The algorithm of the primal-dual interior point method consists of four basic parts: determination of the matrix $G(\mathbf{x}, \mathbf{u})$ or its approximation, solving linear system (11.115), a step length selection, and an update of the barrier parameter μ . The matrix $G(\mathbf{x}, \mathbf{u})$ has form (11.33), so its approximation can be determined in the one of the ways introduced in Remark 11.13.

The linear system (11.115), obtained by determination and subsequent elimination of inactive constraints in the way described in the previous subsection, is solved either directly using the Bunch–Parlett decomposition or iteratively by the conjugate gradient method with the indefinite preconditioner

$$C = \begin{bmatrix} \hat{D} & \hat{A}(\tilde{\mathbf{x}}) \\ \hat{A}^T(\tilde{\mathbf{x}}) & -\hat{M} \end{bmatrix}, \tag{11.123}$$

where \hat{D} is a positive definite diagonal matrix that approximates matrix $\hat{G}(\tilde{\mathbf{x}}, \mathbf{u})$. An iterative process is terminated if residual vectors satisfy conditions (11.122) and

$$\|\mathbf{r}\| \leq \omega \|\tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u})\|, \quad \|\hat{\mathbf{r}}\| \leq \omega \min\{\|\hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \mu \hat{U}^{-1} \tilde{\mathbf{e}}\|, \|\hat{\mathbf{c}}(\tilde{\mathbf{x}}) + \hat{\mathbf{s}}\|\},$$

where $0 < \omega < 1$ is a prescribed precision. The directional derivative $P'(0)$ given by (11.118) should be negative. There are two possibilities how this requirement can be achieved. We either determine the value $\sigma \geq 0$ satisfying inequality (11.121), which implies $P'(0) < 0$ if (11.122) holds (Theorem 11.10), or set $\sigma = 0$ and ignore inequality (11.122). If $P'(0) \geq 0$, we determine a diagonal matrix \tilde{D} with elements

$$\begin{cases} \tilde{D}_{jj} = \underline{\Gamma}, & \text{if } \frac{\|\tilde{\mathbf{g}}\|}{10} |\tilde{G}_{jj}| < \underline{\Gamma}, \\ \tilde{D}_{jj} = \frac{\|\tilde{\mathbf{g}}\|}{10} |\tilde{G}_{jj}|, & \text{if } \underline{\Gamma} \leq \frac{\|\tilde{\mathbf{g}}\|}{10} |\tilde{G}_{jj}| \leq \overline{\Gamma}, \\ \tilde{D}_{jj} = \overline{\Gamma}, & \text{if } \overline{\Gamma} < \frac{\|\tilde{\mathbf{g}}\|}{10} |\tilde{G}_{jj}|, \end{cases} \tag{11.124}$$

for $1 \leq j \leq n + m$, where $\tilde{\mathbf{g}} = \tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u})$ and $0 < \underline{\Gamma} < \overline{\Gamma}$, set $\tilde{G}(\tilde{\mathbf{x}}, \mathbf{u}) = \tilde{D}$ and restart the iterative process by solving new linear system (11.115).

We use functions $\mathbf{s}(\alpha) = [s_{kl}(\alpha)]$, $\mathbf{u}(\alpha) = [u_{kl}(\alpha)]$, where $s_{kl}(\alpha) = s_{kl} + \alpha_{s_{kl}} \Delta s_{kl}$, $u_{kl}(\alpha) = u_{kl} + \alpha_{u_{kl}} \Delta u_{kl}$ and

$$\begin{cases} \alpha_{s_{kl}} = \alpha, & \text{if } \Delta s_{kl} \geq 0, \\ \alpha_{s_{kl}} = \min\{\alpha, -\gamma \frac{s_{kl}}{\Delta s_{kl}}\}, & \text{if } \Delta s_{kl} < 0, \\ \alpha_{u_{kl}} = \alpha, & \text{if } \Delta u_{kl} \geq 0, \\ \alpha_{u_{kl}} = \min\{\alpha, -\gamma \frac{u_{kl}}{\Delta u_{kl}}\}, & \text{if } \Delta u_{kl} < 0 \end{cases}$$

when choosing a step length using the augmented Lagrange function. A parameter $0 < \gamma < 1$ (usually $\gamma = 0.99$) assures the positivity of vectors \mathbf{s}_+ and \mathbf{u}_+ in (11.117). A parameter $\alpha > 0$ is chosen to satisfy the inequality $P(\alpha) - P(0) \leq \varepsilon_1 \alpha P'(0)$, which is possible because $P'(0) < 0$ and a function $P(\alpha)$ is continuous.

After finishing the iterative step, a barrier parameter μ should be updated. There exist several heuristic procedures for this purpose. The following procedure proposed in [28] seems to be very efficient.

Procedure C

Compute the centrality measure

$$\varrho = \frac{\bar{m} \min_{kl} \{s_{kl} u_{kl}\}}{s^T \mathbf{u}},$$

where $\bar{m} = m_1 + \dots + m_m$ and $1 \leq k \leq m, 1 \leq l \leq m_k$. Compute the value

$$\lambda = 0.1 \min \left\{ 0.05 \frac{1-\varrho}{\varrho}, 2 \right\}^3$$

and set $\mu = \lambda s^T \mathbf{u} / \bar{m}$.

The considerations up to now are summarized in Algorithm 11.3 introduced in the Appendix.

11.5 Numerical Experiments

The methods studied in this contribution were tested by using two collections of test problems TEST14 and TEST15 described in [19], which are the parts of the UFO system [24] and can be downloaded from the web-page www.cs.cas.cz/luksan/test.html. Both these collections contain 22 problems with functions $f_k(\mathbf{x})$, $1 \leq k \leq m$, $\mathbf{x} \in \mathbb{R}^n$, where n is an input parameter and $m \geq n$ depends on n (we have used the values $n = 100$ and $n = 1000$ for numerical experiments). Functions $f_k(\mathbf{x})$, $1 \leq k \leq m$, have a sparse structure (the Jacobian matrix of a mapping $\mathbf{f}(\mathbf{x})$ is sparse), so sparse matrix decompositions can be used for solving linear equation systems.

The tested methods, whose results are reported in Tables 11.1, 11.2, 11.3, 11.4, and 11.5 introduced in the Appendix, are denoted by five letters. The first pair of letters gives the problem type: either a classical minimax MX (when a function $F(\mathbf{x})$ has form (11.10) or $F(\mathbf{x}) = \|f(\mathbf{x})\|_\infty$ holds) or a sum of absolute values SA (when $F(\mathbf{x}) = \|f(\mathbf{x})\|_1$ holds). Further two letters specify the method used:

- PI –the primal interior point method (Sect. 11.2),
- SM –the smoothing method (Sect. 11.3),
- DI –the primal-dual interior point method (Sect. 11.4).

The last letter denotes the procedure for updating a barrier parameter μ (procedures A and B are described in Sect. 11.2.4 and procedure C is described in Sect. 11.4.2).

Table 11.1 TEST14 (minimization of maxima) — 22 problems

Method	Newton methods: n=100						Variable metric methods: n=100					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	2232	7265	11575	0.74	4	-	2849	5078	2821	0.32	2	-
MXPI-B	2184	5262	9570	0.60	1	-	1567	2899	1589	0.24	1	-
MXSM-A	3454	11682	21398	1.29	5	-	4444	12505	4465	1.03	-	-
MXSM-B	10241	36891	56399	4.15	3	-	8861	32056	8881	2.21	1	1
MXDI-C	1386	2847	14578	0.90	2	-	2627	5373	2627	0.96	3	-

Method	Newton methods: n=1000						Variable metric methods: n=1000					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	1386	3735	7488	5.58	4	-	3237	12929	3258	5.91	6	-
MXPI-B	3153	6885	12989	9.03	4	-	1522	3287	1544	2.68	5	-
MXSM-A	10284	30783	82334	54.38	7	-	4221	9519	4242	8.00	8	-
MXSM-B	18279	61180	142767	87.76	6	-	13618	54655	13639	45.10	9	1
MXDI-C	3796	6677	48204	49.95	6	-	2371	5548	2371	18.89	3	-

Table 11.2 TEST14 (L_∞ -approximation) — 22 problems

Method	Newton methods: n=100						Variable metric methods: n=100					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	2194	5789	10553	0.67	3	-	2890	5049	2912	0.48	1	-
MXPI-B	6767	17901	39544	3.79	4	-	1764	3914	1786	0.37	2	-
MXSM-A	3500	9926	23568	1.79	7	-	8455	23644	8476	4.69	4	-
MXSM-B	15858	48313	92486	8.33	5	-	9546	34376	9566	2.59	9	1
MXDI-C	1371	2901	11580	1.12	8	-	2467	5130	2467	1.59	3	-

Method	Newton methods: n=1000						Variable metric methods: n=1000					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	4110	14633	20299	18.89	4	-	1549	2636	1571	2.51	3	-
MXPI-B	6711	31618	29939	30.73	7	-	1992	6403	2013	4.96	4	-
MXSM-A	9994	24333	88481	67.45	11	-	6164	15545	6185	29.37	8	-
MXSM-B	23948	84127	182604	149.63	8	-	24027	92477	24048	132.08	8	1
MXDI-C	3528	9084	26206	49.78	12	-	1932	2845	1932	18.73	5	-

The columns of all tables correspond to two classes of methods. The Newton methods use approximations of the Hessian matrices of the Lagrange function obtained by gradient differences [4] and variable metric methods update approximations of the Hessian matrices of the partial functions by the methods belonging to the Broyden family [12] (Remark 11.13).

Table 11.3 TEST15 (L_∞ -approximation) — 22 problems

Method	Newton methods: n=100						Variable metric methods: n=100					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	15525	20272	55506	4.41	1	-	6497	8204	6518	1.37	3	-
MXPI-B	7483	17999	27934	3.27	5	-	1764	7598	2488	0.74	2	-
MXSM-A	17574	29780	105531	11.09	4	-	9879	15305	9900	5.95	-	-
MXSM-B	13446	29249	81938	6.80	9	1	9546	34376	9566	2.59	3	-
MXDI-C	980	1402	7356	0.79	1	-	1179	1837	1179	1.06	2	-
Method	Newton methods: n=1000						Variable metric methods: n=1000					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
MXPI-A	10325	15139	32422	39.30	6	-	6484	9904	6502	13.77	2	-
MXPI-B	14836	30724	46864	68.70	10	-	7388	15875	7409	19.98	3	-
MXSM-A	11722	24882	69643	61.65	10	-	6659	12824	6681	41.55	8	-
MXSM-B	13994	31404	86335	78.45	9	1	15125	25984	15147	61.62	10	-
MXDI-C	1408	2406	10121	15.63	6	-	2228	3505	2228	35.13	10	-

Table 11.4 TEST14 (L_1 -approximation) — 22 problems

Method	Newton methods: n=100						Variable metric methods: n=100					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
SAPI-A	1647	5545	8795	0.63	5	-	12265	23579	12287	1.37	2	1
SAPI-B	1957	7779	10121	0.67	6	-	4695	6217	10608	0.67	3	-
SASM-A	1677	4505	16079	0.74	3	-	20025	27369	20047	2.83	4	-
SASM-B	2389	8085	23366	1.18	4	-	5656	11637	5678	1.02	2	-
SADI-C	4704	13012	33937	4.16	7	1	6547	7012	6547	9.18	8	-
Method	Newton methods: n=1000						Variable metric methods: n=1000					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
SAPI-A	7592	19621	46100	15.39	4	-	22277	36610	22298	19.09	7	1
SAPI-B	9067	35463	56292	19.14	6	-	16650	35262	16672	14.47	6	1
SASM-A	5696	13534	41347	15.28	4	-	20020	30736	20042	23.05	5	1
SASM-B	8517	30736	57878	23.60	6	-	18664	28886	18686	18.65	5	1
SADI-C	6758	11011	47960	94.78	11	1	13123	14610	13124	295.46	8	2

The tables contain total numbers of iterations NIT, function evaluations NFV, gradient evaluations NFG, and also the total computational time, the number of problems with the value $\bar{\Delta}$ decreased and the number of failures (the number of unsolved problems). The decrease of the maximum step length $\bar{\Delta}$ is used for three reasons. First, too large steps may lead to overflows if arguments of functions (roots,

Table 11.5 TEST15 (L_1 -approximation) — 22 problems

Method	Newton methods: n=100						Variable metric methods: n=100					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
SAPI-A	15760	21846	58082	4.24	8	-	39469	58157	39486	6.28	4	1
SAPI-B	4592	17050	17778	1.46	5	-	5932	25035	5952	1.48	6	1
SASM-A	10098	14801	610511	3.54	5	-	9162	28421	9184	3.65	6	1
SASM-B	4528	14477	290379	2.94	8	-	3507	9036	3528	1.27	6	-
SADI-C	877	1373	6026	0.84	3	-	15528	15712	15529	14.49	5	1

Method	Newton methods: n=1000						Variable metric methods: n=1000					
	NIT	NFV	NFG	Time	Δ	Fail	NIT	NFV	NFG	Time	Δ	Fail
SAPI-A	18519	39319	70951	61.04	5	-	27308	44808	27327	36.64	4	1
SAPI-B	12405	57969	43189	55.06	7	-	12712	32179	12731	21.48	8	1
SASM-A	19317	32966	113121	62.65	8	-	22264	42908	22284	62.46	7	1
SASM-B	14331	33572	86739	57.56	6	-	12898	42479	12919	47.05	7	1
SADI-C	2093	3681	12616	20.01	3	1	23957	28000	23960	186.92	5	3

logarithms, exponentials) lie outside of their definition domain. Second, the change of Δ can affect the finding of a suitable (usually global) minimum. Finally, it can prevent from achieving a domain in which the objective function has bad behavior leading to a loss of convergence. The number of times the step length has decreased is in some sense a symptom of robustness (a lower number corresponds to higher robustness).

Several conclusions can be done from the results stated in these tables.

- The smoothing methods are less efficient than the primal interior point methods. For testing the smoothing methods, we had to use the value $\underline{\mu} = 10^{-6}$, while the primal interior methods work well with the smaller value $\underline{\mu} = 10^{-8}$, which gives more precise results.
- The primal-dual interior point methods are slower than the primal interior point methods, especially for the reason that system of equations (11.115) is indefinite, so we cannot use the Choleski (or the Gill–Murray [10]) decomposition. If the matrix of linear system (11.115) is large and sparse, we can use the Bunch–Parlett decomposition [6]. In this case, a large fill-in of new nonzero elements (and thus to overflow of the operational memory or large extension of the computational time) may appear. In this case, we can also use the iterative conjugate gradient method with an indefinite preconditioner [18], however, the ill-conditioned systems can require a large number of iterations and thus a large computational time.
- It cannot be uniquely decided whether Procedure A is better than Procedure B. The Newton methods usually work better with Procedure A while the variable metric methods are more efficient with Procedure B.

- The variable metric methods are usually faster because it is not necessary to determine the elements of the Hessian matrix of the Lagrange function by gradient differences. The Newton methods seem to be more robust (especially in case of L_1 -approximation).

Appendix

Algorithm 11.1: Primal interior point method

Data: A tolerance for the gradient norm of the Lagrange function $\underline{\varepsilon} > 0$. A precision for determination of a minimax vector $\underline{\delta} > 0$. Bounds for a barrier parameter $0 < \underline{\mu} < \bar{\mu}$. Coefficients for decrease of a barrier parameter $0 < \lambda < 1, \sigma > 1$ (or $0 < \vartheta < 1$). A tolerance for a uniform descent $\varepsilon_0 > 0$. A tolerance for a step length selection $\varepsilon_1 > 0$. A maximum step length $\bar{\Delta} > 0$.

Input. A sparsity pattern of the matrix $A(\mathbf{x}) = [A_1(\mathbf{x}), \dots, A_m(\mathbf{x})]$. A starting point $\mathbf{x} \in \mathbb{R}^n$.

Step 1. (Initiation) Choose $\mu \leq \bar{\mu}$. Determine a sparse structure of the matrix $W = W(\mathbf{x}; \mu)$ from the sparse structure of the matrix $A(\mathbf{x})$ and perform a symbolic decomposition of the matrix W (described in [2, Section 1.7.4]). Compute values $f_{kl}(\mathbf{x}), 1 \leq k \leq m, 1 \leq l \leq m_k$, values $F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x}), 1 \leq k \leq m$, and the value of objective function (11.4). Set $r = 0$ (restart indicator).

Step 2. (Termination) Solve nonlinear equations (11.44) with precision $\underline{\delta}$ to obtain a minimax vector $\mathbf{z}(\mathbf{x}; \mu)$ and a vector of Lagrange multipliers $\mathbf{u}(\mathbf{x}; \mu)$. Determine a matrix $A = A(\mathbf{x})$ and a vector $\mathbf{g} = \mathbf{g}(\mathbf{x}; \mu) = A(\mathbf{x})\mathbf{u}(\mathbf{x}; \mu)$. If $\mu \leq \underline{\mu}$ and $\|\mathbf{g}\| \leq \underline{\varepsilon}$, terminate the computation.

Step 3. (Hessian matrix approximation) Set $G = G(\mathbf{x}; \mu)$ or compute an approximation G of the Hessian matrix $G(\mathbf{x}; \mu)$ using gradient differences or using quasi-Newton updates (Remark 11.13).

Step 4. (Direction determination) Determine a matrix $\nabla^2 \hat{B}(\mathbf{x}; \mu)$ by (11.48) and a vector $\Delta \mathbf{x}$ by solving Eq. (11.49) with the right-hand side defined by (11.47).

Step 5. (Restart) If $r = 0$ and (11.54) does not hold, set $G = I, r = 1$ and go to Step 4. If $r = 1$ and (11.54) does not hold, set $\Delta \mathbf{x} = -\mathbf{g}$. Set $r = 0$.

Step 6. (Step length selection) Determine a step length $\alpha > 0$ satisfying inequalities (11.55) (for a barrier function $\hat{B}(\mathbf{x}; \mu)$ defined by (11.46)) and $\alpha \leq \bar{\Delta} / \|\Delta \mathbf{x}\|$. Note that nonlinear equations (11.44) are solved at the point $\mathbf{x} + \alpha \Delta \mathbf{x}$. Set $\mathbf{x} := \mathbf{x} + \alpha \Delta \mathbf{x}$. Compute values $f_{kl}(\mathbf{x}), 1 \leq k \leq m, 1 \leq l \leq m_k$, values $F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x}), 1 \leq k \leq m$, and the value of objective function (11.4).

Step 7. (Barrier parameter update) Determine a new value of a barrier parameter $\mu \geq \underline{\mu}$ using Procedure A or Procedure B. Go to Step 2.

The values $\underline{\varepsilon} = 10^{-6}$, $\underline{\delta} = 10^{-6}$, $\underline{\mu} = 10^{-8}$, $\overline{\mu} = 1$, $\lambda = 0.85$, $\sigma = 100$, $\vartheta = 0.1$, $\varepsilon_0 = 10^{-8}$, $\varepsilon_1 = 10^{-4}$, and $\overline{\Delta} = 1000$ were used in our numerical experiments.

Algorithm 11.2: Smoothing method

Data: A tolerance for the gradient norm of the smoothing function $\underline{\varepsilon} > 0$.

Bounds for a smoothing parameter $0 < \underline{\mu} < \overline{\mu}$. Coefficients for decrease of a smoothing parameter $0 < \lambda < 1$, $\sigma > 1$ (or $0 < \vartheta < 1$). A tolerance for a uniform descent $\varepsilon_0 > 0$. A tolerance for a step length selection $\varepsilon_1 > 0$. A maximum step length $\overline{\Delta} > 0$.

Input. A sparsity pattern of the matrix $A(\mathbf{x}) = [A_1(\mathbf{x}), \dots, A_m(\mathbf{x})]$. A starting point $\mathbf{x} \in \mathbb{R}^n$.

Step 1. (*Initiation*) Choose $\mu \leq \overline{\mu}$. Determine a sparse structure of the matrix $W = W(\mathbf{x}; \mu)$ from the sparse structure of the matrix $A(\mathbf{x})$ and perform a symbolic decomposition of the matrix W (described in [2, Section 1.7.4]). Compute values $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, values $F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, and the value of objective function (11.4). Set $r = 0$ (restart indicator).

Step 2. (*Termination*) Determine a vector of smoothing multipliers $\mathbf{u}(\mathbf{x}; \mu)$ by (11.87). Determine a matrix $A = A(\mathbf{x})$ and a vector $\mathbf{g} = \mathbf{g}(\mathbf{x}; \mu) = A(\mathbf{x})\mathbf{u}(\mathbf{x}; \mu)$. If $\mu \leq \underline{\mu}$ and $\|\mathbf{g}\| \leq \underline{\varepsilon}$, terminate the computation.

Step 3. (*Hessian matrix approximation*) Set $G = G(\mathbf{x}; \mu)$ or compute an approximation G of the Hessian matrix $G(\mathbf{x}; \mu)$ using gradient differences or using quasi-Newton updates (Remark 11.13).

Step 4. (*Direction determination*) Determine the matrix W by (11.94) and the vector $\Delta\mathbf{x}$ by (11.93) using the Gill–Murray decomposition of the matrix W .

Step 5. (*Restart*) If $r = 0$ and (11.54) does not hold, set $G = I$, $r = 1$ and go to Step 4. If $r = 1$ and (11.54) does not hold, set $\Delta\mathbf{x} = -\mathbf{g}$. Set $r = 0$.

Step 6. (*Step length selection*) Determine a step length $\alpha > 0$ satisfying inequalities (11.55) (for a smoothing function $S(\mathbf{x}; \mu)$) and $\alpha = \overline{\Delta}/\|\Delta\mathbf{x}\|$. Set $\mathbf{x} := \mathbf{x} + \alpha\Delta\mathbf{x}$. Compute values $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, values $F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, and the value of the objective function (11.4).

Step 7. (*Smoothing parameter update*) Determine a new value of the smoothing parameter $\mu \geq \underline{\mu}$ using Procedure A or Procedure B. Go to Step 2.

The values $\underline{\varepsilon} = 10^{-6}$, $\underline{\mu} = 10^{-6}$, $\overline{\mu} = 1$, $\lambda = 0.85$, $\sigma = 100$, $\vartheta = 0.1$, $\varepsilon_0 = 10^{-8}$, $\varepsilon_1 = 10^{-4}$, and $\overline{\Delta} = 1000$ were used in our numerical experiments.

Algorithm 11.3: Primal-dual interior point method

Data: A tolerance for the gradient norm $\underline{\varepsilon} > 0$. A parameter for determination of active constraints $\tilde{\varepsilon} > 0$. A parameter for initiation of slack variables and Lagrange multipliers $\delta > 0$. An initial value of the barrier parameter $\bar{\mu} > 0$. A precision for the direction determination $0 \leq \omega < 1$. A parameter for the step length selection $0 < \gamma < 1$. A tolerance for the step length selection $\varepsilon_1 > 0$. Maximum step length $\bar{\Delta} > 0$.

Input. A sparsity pattern of the matrix $A(\mathbf{x}) = [A_1(\mathbf{x}), \dots, A_m(\mathbf{x})]$. A starting point $\mathbf{x} \in \mathbb{R}^n$.

Step 1. (Initialization) Compute values $f_{kl}(\mathbf{x})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and set $F_k(\mathbf{x}) = \max_{1 \leq l \leq m_k} f_{kl}(\mathbf{x})$, $z_k = F_k(\mathbf{x}) + \delta$, $1 \leq k \leq m$. Compute values $c_{kl}(\tilde{\mathbf{x}}) = f_{kl}(\tilde{\mathbf{x}}) - z_k$, and set $s_{kl} = -c_{kl}(\tilde{\mathbf{x}})$, $u_{kl} = \delta$. Set $\mu = \bar{\mu}$ and compute the value of the barrier function $\tilde{B}_\mu(\tilde{\mathbf{x}}, \mathbf{s})$.

Step 2. (Termination) Determine a matrix $\tilde{A}(\tilde{\mathbf{x}})$ and a vector $\tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u}) = \tilde{A}(\tilde{\mathbf{x}})\mathbf{u}$ by (11.111). If the KKT conditions $\|\tilde{\mathbf{g}}(\tilde{\mathbf{x}}, \mathbf{u})\| \leq \underline{\varepsilon}$, $\|\mathbf{c}(\tilde{\mathbf{x}}) + \mathbf{s}\| \leq \underline{\varepsilon}$, and $\mathbf{s}^T \mathbf{u} \leq \underline{\varepsilon}$ are satisfied, terminate the computation.

Step 3. (Hessian matrix approximation) Set $G = G(\mathbf{x}, \mathbf{u})$ or compute an approximation G of the Hessian matrix $G(\mathbf{x}, \mathbf{u})$ using gradient differences or utilizing quasi-Newton updates (Remark 11.13). Determine a parameter $\sigma \geq 0$ by (11.121) or set $\sigma = 0$. Split the constraints into active if $\hat{s}_{kl} \leq \tilde{\varepsilon} \hat{u}_{kl}$ and inactive if $\check{s}_{kl} > \tilde{\varepsilon} \check{u}_{kl}$.

Step 4. (Direction determination) Determine the matrix $\tilde{G} = \tilde{G}(\tilde{\mathbf{x}}, \mathbf{u})$ by (11.111) (where the Hessian matrix $G(\mathbf{x}, \mathbf{u})$ is replaced with its approximation G). Determine vectors $\Delta\tilde{\mathbf{x}}$ and $\Delta\tilde{\mathbf{u}}$ by solving linear system (11.115), a vector $\Delta\tilde{\mathbf{u}}$ by (11.114), and a vector $\Delta\mathbf{s}$ by (11.116). Linear system (11.115) is solved either directly using the Bunch–Parlett decomposition (we carry out both the symbolic and the numeric decompositions in this step) or iteratively by the conjugate gradient method with indefinite preconditioner (11.123). Compute the derivative of the augmented Lagrange function by formula (11.120).

Step 5. (Restart) If $P'(0) \geq 0$, determine a diagonal matrix \tilde{D} by (11.124), set $\tilde{G} = \tilde{D}$, $\sigma = 0$, and go to Step 4.

Step 6. (Step length selection) Determine a step length parameter $\alpha > 0$ satisfying inequalities $P(\alpha) - P(0) \leq \varepsilon_1 \alpha P'(0)$ and $\alpha \leq \bar{\Delta} / \|\Delta\mathbf{x}\|$. Determine new vectors $\tilde{\mathbf{x}} := \tilde{\mathbf{x}} + \alpha \Delta\tilde{\mathbf{x}}$, $\mathbf{s} := \mathbf{s}(\alpha)$, $\mathbf{u} := \mathbf{u}(\alpha)$ by (11.117). Compute values $f_{kl}(\tilde{\mathbf{x}})$, $1 \leq k \leq m$, $1 \leq l \leq m_k$, and set $c_{kl}(\tilde{\mathbf{x}}) = f_{kl}(\tilde{\mathbf{x}}) - z_k$, $1 \leq k \leq m$, $1 \leq l \leq m_k$. Compute the value of the barrier function $\tilde{B}_\mu(\tilde{\mathbf{x}}, \mathbf{s})$.

Step 7. (Barrier parameter update) Determine a new value of the barrier parameter $\mu \geq \underline{\mu}$ using Procedure C. Go to Step 2.

The values $\underline{\varepsilon} = 10^{-6}$, $\tilde{\varepsilon} = 0.1$, $\delta = 0.1$, $\omega = 0.9$, $\gamma = 0.99$, $\bar{\mu} = 1$, $\varepsilon_1 = 10^{-4}$, and $\bar{\Delta} = 1000$ were used in our numerical experiments.

References

1. Benzi, M., Golub, G.H., Liesen, J.: Numerical solution of saddle point problems. *Acta Numer.* **14**, 1–137 (2005)
2. Björck, A.: *Numerical Methods in Matrix Computations*. Springer, New York (2015)
3. Bunch, J.R., Parlett, B.N.: Direct methods for solving symmetric indefinite systems of linear equations. *SIAM J. Numer. Anal.* **8**, 639–655 (1971)
4. Coleman, T.F., Moré, J.J.: Estimation of sparse Hessian matrices and graph coloring problems. *Math. Program.* **28**, 243–270 (1984)
5. Di Pillo, G., Grippo, L., Lucidi, S.: Smooth Transformation of the generalized minimax problem. *J. Optim. Theory Appl.* **95**, 1–24 (1997)
6. Duff, I.S., Gould, N.I.M., Reid, J.K., Turner, K.: The factorization of sparse symmetric indefinite matrices. *IMA J. Numer. Anal.* **11**, 181–204 (1991)
7. Fiedler, M.: *Special Matrices and Their Applications in Numerical Mathematics*. Dover, New York (2008)
8. Fletcher, R.: *Practical Methods of Optimization*. Wiley, New York (1987)
9. Fletcher, R., Leyffer, S.: Nonlinear programming without a penalty function. *Math. Program.* **91**, 239–269 (2002)
10. Gill, P.E., Murray, W.: Newton type methods for unconstrained and linearly constrained optimization. *Math. Program.* **7**, 311–350 (1974)
11. Gould, N.I.M., Toint, P.L.: Nonlinear programming without a penalty function or a filter. *Math. Program.* **122**, 155–196 (2010)
12. Griewank, A., Toint, P.L.: Partitioned variable metric updates for large-scale structured optimization problems. *Numer. Math.* **39**, 119–137 (1982)
13. Griewank, A., Walther, A.: *Evaluating Derivatives*. SIAM, Philadelphia (2008)
14. Le, D.: Three new rapidly convergent algorithms for finding a zero of a function. *SIAM J. Sci. Stat. Comput.* **6**, 193–208 (1985)
15. Le, D.: An efficient derivative-free method for solving nonlinear equations. *ACM Trans. Math. Softw.* **11**, 250–262 (1985)
16. Liu, X., Yuan, Y.: A sequential quadratic programming method without a penalty function or a filter for nonlinear equality constrained optimization. *SIAM J. Optim.* **21**, 545–571 (2011)
17. Lukšan, L., Spedicato, E.: Variable metric methods for unconstrained optimization and nonlinear least squares. *J. Comput. Appl. Math.* **124**, 61–93 (2000)
18. Lukšan, L., Vlček, J.: Indefinitely preconditioned inexact Newton method for large sparse equality constrained nonlinear programming problems. *Numer. Linear Algebra Appl.* **5**, 219–247 (1998)
19. Lukšan, L., Vlček, J.: Sparse and partially separable test problems for unconstrained and equality constrained optimization. Technical Report V-767, Prague, ICS AS CR (1998)
20. Lukšan, L., Matonoha, C., Vlček, J.: Interior-point method for non-linear non-convex optimization. *Numer. Linear Algebra Appl.* **11**, 431–453 (2004)
21. Lukšan, L., Matonoha, C., Vlček, J.: Interior-point method for large-scale nonlinear programming. *Optim. Methods Softw.* **20**, 569–582 (2005)
22. Lukšan, L., Matonoha, C., Vlček, J.: Primal interior-point method for large sparse minimax optimization. *Kybernetika* **45**, 841–864 (2009)
23. Lukšan, L., Matonoha, C., Vlček, J.: Primal interior-point method for minimization of generalized minimax functions. *Kybernetika* **46**, 697–721 (2010)
24. Lukšan, L., Tůma, M., Matonoha, C., Vlček, J., Ramešová, N., Šiška, M., Hartman, J.: UFO 2017. Interactive System for Universal Functional Optimization. Technical Report V-1252, Prague, ICS AS CR (2017)
25. Mäkelä, M.M., Neittaanmäki, P.: *Nonsmooth Optimization*. World Scientific, London (1992)
26. Nocedal, J., Wright, S.J.: *Numerical Optimization*. Springer, New York (2006)

27. Tůma, M.: A note on direct methods for approximations of sparse Hessian matrices. *Appl. Math.* **33**, 171–176 (1988)
28. Vanderbei, J., Shanno, D.F.: An interior point algorithm for nonconvex nonlinear programming. *Comput. Optim. Appl.* **13**, 231–252 (1999)
29. Wächter, A., Biegler, L.: Line search filter methods for nonlinear programming. Motivation and global convergence. *SIAM J. Comput.* **16**, 1–31 (2005)