# SPECTRAL BUNDLE METHODS FOR NON-CONVEX MAXIMUM EIGENVALUE FUNCTIONS. PART 1: FIRST-ORDER METHODS

D. NOLL [*] AND P. APKARIAN [†]

**Abstract.** Many challenging problems in automatic control may be cast as optimization programs subject to matrix inequality constraints. Here we investigate an approach which converts such problems into non-convex eigenvalue optimization programs and makes them amenable to nonsmooth analysis techniques like bundle or cutting plane methods. We prove global convergence of a first order bundle method for programs with non-convex maximum eigenvalue functions.

**Key words.** Bilinear matrix inequality (BMI), linear matrix inequality (LMI), eigenvalue optimization, first-order spectral bundle method, $\epsilon$-subgradient.

**1. Introduction.** The importance of linear matrix inequalities (LMIs) and bilinear matrix inequalities (BMIs) for applications in automatic control has been recognized during the past decade. Semidefinite programming (SDP) used to solve LMI problems has found a widespread interest due to its large spectrum of applications [8]. But many challenging engineering design problems lead to BMI feasibility or optimization programs no longer amenable to convex methods. Some prominent BMI problems are parametric robust feedback control [60, 61], static and reduced-order controller design [1, 16], design of structured controllers [15], decentralized synthesis, or synthesis with finite precision controllers [70]. These problems are in fact known to be NP-hard (see e.g. [48]).

Due to their significance for industrial applications, many solution strategies for BMI problems have been proposed. The most ambitious ones use ideas from global optimization such as branch-and-bound methods [21, 9, 23] or concave programming [5, 6] in order to address the presence of multiple local minima. On the other hand, in many situations, semidefinite programming relaxations, or heuristic approaches like coordinate descent schemes or alternating techniques (analysis versus synthesis) have been used with considerable success; see [11, 69, 29, 22, 32].

The general nature of BMI problems, which include for instance quadratic constraint quadratic (QCQP) programming, all polynomial problems and mixed binary programming, makes it evident that a general strategy may hardly be expected. We have observed that BMIs in control applications may usually be solved by local methods. This is significant, since global methods have a prohibitive computational load and are therefore of very limited applicability. We have contributed several local nonlinear programming approaches, which are capable to deal with matrix inequality constraints [1, 2, 16, 17, 50]; see [31] for another approach. Here we investigate a strategy which converts BMI problems into non-convex eigenvalue optimization programs, which are then solved by non-smooth analysis tools.

**1.1. BMIs and eigenvalue programs.** We consider affine $\mathcal{A} : \mathbb{R}^n \to \mathbb{S}^m$ and bilinear operators $\mathcal{B} : \mathbb{R}^n \to \mathbb{S}^m$ into the space $\mathbb{S}^m$ of symmetric $m \times m$ matrices,

$$(1.1) \qquad \mathcal{A}(x) = A_0 + \sum_{i=1}^{n} A_i x_i, \quad \mathcal{B}(x) = \mathcal{A}(x) + \sum_{1 \leq i < j \leq n} B_{ij} x_i x_j.$$

---

[*]Université Paul Sabatier, Institut de Mathématiques de Toulouse (IMT), 118, route de Narbonne, 31062 Toulouse, France

[†]CERT-ONERA, Control System Department, 2, avenue Edouard Bélin, 31055 Toulouse, France

Then a BMI-optimization program is of the form

$$(1.2) \qquad \begin{array}{ll} \text{minimize} & c^\top x,\ x \in \mathbb{R}^n \\ \text{subject to} & \mathcal{B}(x) \preceq 0 \end{array}$$

where $\preceq$ means negative semidefinite. Semidefinite programming is a special case where an affine operator $\mathcal{A}$ is used. The simpler BMI-feasibility problem seeks $x \in \mathbb{R}^n$ satisfying $\mathcal{B}(x) \preceq 0$. This is a special case of (1.2) if solved as $\min\{t : \mathcal{B}(x) \preceq tI_m\}$.

These problems are clearly related to eigenvalue optimization. We consider the unconstrained

$$(1.3) \qquad \text{minimize} \quad \lambda_1\left(\mathcal{B}(x)\right),\ x \in \mathbb{R}^n,$$

and the constrained eigenvalue optimization program

$$(1.4) \qquad \begin{array}{ll} \text{minimize} & c^\top x,\ x \in \mathbb{R}^n \\ \text{subject to} & \lambda_1\left(\mathcal{B}(x)\right) \leq 0. \end{array}$$

Here $\lambda_1 : \mathbb{S}^m \to \mathbb{R}$ is the maximum eigenvalue function, which is convex but nonsmooth in general. The non-convexity of (1.3) and (1.4) is induced by the operator $\mathcal{B}$. Clearly (1.2) is equivalent to (1.4), while the BMI-feasibility problem may be solved using (1.3). Programs of the form (1.4) may be transformed into (1.3) via exact penalization, even though it may be preferable to use the structure of (1.4) explicitly.

Eigenvalue optimization has an interest of its own even outside control applications. Much pioneering work has been contributed by M. Overton in a series of papers [51, 52, 53, 54, 55] beginning in the 1980s, where Newton type methods for (1.3) are considered. Further to be mentioned among the earliest contributions are J. Cullum *et al.* [13], R. Fletcher [18], A. Shapiro [63, 64] and A. Shapiro and M.K.H. Fan [65].

Bundle methods have been invented by C. Lemaréchal [40] and Wolfe [71] and developed mostly in the 1980s by numerous contributions, in particular from Lemaréchal [41, 43] and K. Kiwiel (see e.g. [34, 37, 36]). A survey is [42], see also [30, 34, 62]. Bundle methods have recently been revived in the context of semidefinite programming; see [58, 44, 26, 46, 47, 39, 59]. Early contributions to non-convex programs are Kiwiel [37, 35, 38], a recent non-convex bundle algorithm is presented in Fuduli *et al.* [20], a convex filter bundle method is given in [19].

**1.2. Purpose.** As a rule, bundle methods maintain a finite set of affine estimates of the objective function $f$ at the current iterate. This set is used to predict a descent step. This trial step is either accepted (serious step) if the actual descent is sufficient, or it is used to improve the local model (null step) by adding another affine approximation of $f$. It has been observed in [13] (and used in [57, 58, 26, 44] and also [24, 27, 28, 25]) that maximum eigenvalue functions $f$ allow for specific bundling strategies, where infinite sets of affine estimates of $f$ are manageable. Here we develop a similar strategy for non-convex eigenvalue functions. Compared to the convex case, the main difficulty is that approximate subgradients of $f$ do no longer provide global information. This complicates the analysis of the bundling procedure. In section 3 we develop a first-order algorithm for (1.3), which extends the approach of [13] and [58] to the non-convex case. We then solve (1.4) using an improvement function as proposed in [34].

**2. Preparation and preliminary results.** We recall some general notions and definitions and prepare the setting for the analysis of eigenvalue functions.

**2.1. General definitions.** Along with the operators $\mathcal{A}$, $\mathcal{B}$ of (1.1) we consider more general matrix-valued operators $\mathcal{F} : \mathbb{R}^n \to \mathbb{S}^m$ of class $\mathcal{C}^2$. For such $\mathcal{F}$ let $\mathcal{F}'(x) : \mathbb{R}^n \to \mathbb{S}^m$ denote its derivative, whose adjoint $\mathcal{F}'(x)^*$ maps $\mathbb{S}^m$ into $\mathbb{R}^n$. For affine $\mathcal{A}$ we write $A := \mathcal{A}'$, which is the linear part of $\mathcal{A}$, defined as $Ax = \sum_{i=1}^n A_i x_i$. Its adjoint is $A^* : \mathbb{S}^m \to \mathbb{R}^n$, defined as $A^* Z = (A_1 \bullet Z, \ldots, A_n \bullet Z)$. The scalar product in $\mathbb{S}^m$ is $X \bullet Z = \mathrm{tr}(XZ)$. The second derivative $\mathcal{F}''(x)$ of $\mathcal{F}$ is a linear operator $\mathbb{R}^n \to L(\mathbb{R}^n, \mathbb{S}^m)$. For bilinear $\mathcal{B}$ it is independent of $x$ and defined via the tensor $[d, \mathcal{B}'' d] = 2 \sum_{1 \le i < j \le n} B_{ij} d_i d_j \in \mathbb{S}^m$. With these preparations we are ready to consider the class of (generally non-convex) maximum eigenvalue functions $f$ of the form $f = \lambda_1 \circ \mathcal{F}$. Only in the case of an affine operator $\mathcal{A}$ is $f = \lambda_1 \circ \mathcal{A}$ convex and has been analyzed by many authors; see in particular [13, 57, 58].

We use notions from convexity and non-smooth analysis as in [30] or [12]. Given $\epsilon \ge 0$ and a convex function $\phi$ on some Euclidean space $E$, following [30, ch. XI] the $\epsilon$-subdifferential of $\phi$ at $x$ is

$$\partial_\epsilon \phi(x) = \{y \in E^* : \langle y, h \rangle \le \phi(x + h) - \phi(x) - \epsilon \text{ for all } h \in E\}.$$

For $\epsilon = 0$ this is the usual subdifferential in convex analysis. The $\epsilon$-subdifferential gives rise to the $\epsilon$-directional derivative, (cf. [30, ch. XI])

$$\phi_\epsilon'(x; d) = \max\{\langle y, d \rangle : y \in \partial_\epsilon \phi(x)\} = \inf_{t > 0} \frac{\phi(x + td) - \phi(x) - \epsilon}{t},$$

where again $\epsilon = 0$ reproduces the usual directional derivative $\phi'(x; d)$ of convex analysis. Notice that the $\epsilon$-subdifferential of the maximum eigenvalue function $\lambda_1$ is (cf. [30, Example XI.1.2.5]):

(2.1) $$\partial_\epsilon \lambda_1(X) = \{G \in \mathbb{S}^m : G \succeq 0, \mathrm{tr}(G) = 1, G \bullet X \ge \lambda_1(X) - \epsilon\},$$

which is a consequence of the well-known fact that $\lambda_1$ is the support function of the set $\mathcal{C}_m = \{X \in \mathbb{S}^m : X \succeq 0, \mathrm{tr}(X) = 1\}$.

**2.2. Concepts relating to $f = \lambda_1 \circ \mathcal{F}$.** Let us introduce a notation which we will use systematically. Given $x, d \in \mathbb{R}^n$, let $X = \mathcal{F}(x) \in \mathbb{S}^m$ and $D = \mathcal{F}'(x)d \in \mathbb{S}^m$. If $x_k$, $d_k$ arise in an algorithm, we use $X_k$ and $D_k$ accordingly.

Let us now extend $\epsilon$-subgradients and $\epsilon$-directional derivatives to functions $f = \lambda_1 \circ \mathcal{F}$. Let

$$\partial_\epsilon f(x) = \mathcal{F}'(x)^*[\partial_\epsilon \lambda_1(\mathcal{F}(x))] = \mathcal{F}'(x)^*[\partial_\epsilon \lambda_1(X)].$$

Then the $\epsilon$-directional derivative is defined as

$$\begin{aligned}
f_\epsilon'(x; d) &= \max\{g^\top d : g \in \partial_\epsilon f(x)\} \\
&= \max\{G \bullet D : G \in \partial_\epsilon \lambda_1(X)\} \\
&= (\lambda_1)_\epsilon'(X; D).
\end{aligned}$$

It is well-known that the $\epsilon$-directional derivative is difficult to compute in practice. Following Cullum et al. [13] and Oustry [58], we consider the so-called $\epsilon$-enlargement $\delta_\epsilon \lambda_1(X)$ of the subdifferential of the maximum eigenvalue function, which is somewhere in between the purely local subdifferential $\partial \lambda_1(X)$ and the global $\partial_\epsilon \lambda_1(X)$ and turns out a good estimate of the latter.

Following [13] and in particular [58, Def. 1], fix $\epsilon > 0$ and let $r(\epsilon)$ be the largest index $i$ such that $\lambda_i(X) > \lambda_1(X) - \epsilon$, called the $\epsilon$-multiplicity of $\lambda_1(X)$. Notice that

$r(\epsilon)$ is always at the end of a block of equal eigenvalues, that is, $\lambda_1(X) \geq \cdots \geq \lambda_{t-1}(X) > \lambda_t(X) = \cdots = \lambda_{r(\epsilon)}(X) > \lambda_{r(\epsilon)+1}(X) \geq \cdots \geq \lambda_m(X)$, where of course $t = 1$ and $t = r(\epsilon)$ are admitted. We therefore define the spectral separation of $\epsilon$ as

$$(2.2) \qquad \Delta_\epsilon(X) = \lambda_{r(\epsilon)}(X) - \lambda_{r(\epsilon)+1}(X) > 0.$$

Let $Q_\epsilon$ be a $m \times r(\epsilon)$ matrix whose columns form an orthonormal basis of the invariant subspace of $X$ spanned by the eigenvectors of the first $r(\epsilon)$ eigenvalues of $X$. Now define

$$\delta_\epsilon \lambda_1(X) = \{G : G = Q_\epsilon Y Q_\epsilon^\top, \ Y \succeq 0, \ \mathrm{tr}(Y) = 1, \ Y \in \mathbb{S}^{r(\epsilon)}\},$$

then $\delta_\epsilon \lambda_1(X) \subset \partial_\epsilon \lambda_1(X)$. We extend this concept to the class of functions $f = \lambda_1 \circ \mathcal{F}$ by setting

$$\delta_\epsilon f(x) = \mathcal{F}'(x)^*[\delta_\epsilon \lambda_1(X)],$$

so $\partial f(x) \subset \delta_\epsilon f(x) \subset \partial_\epsilon f(x)$, and the $\epsilon$-enlarged subdifferential $\delta_\epsilon f(x)$ may be considered as an inner approximation of the $\epsilon$-subdifferential $\partial_\epsilon f(x)$.

There is a natural analogue of the $\epsilon$-directional derivative based on the new set $\delta_\epsilon f(x)$. Indeed, following [13, 58], define the $\epsilon$-enlarged directional derivative of $\lambda_1$ by

$$(\tilde{\lambda}_1)'_\epsilon(X; D) = \max\{G \bullet D : D \in \delta_\epsilon \lambda_1(X)\} = \sigma_{\delta_\epsilon \lambda_1(X)}(D),$$

where $\sigma_K$ denotes the support function of a convex set $K$. Extend this to functions $f = \lambda_1 \circ \mathcal{F}$ by setting

$$\begin{aligned}
\tilde{f}'_\epsilon(x; d) &= \max\{g^\top d : g \in \delta_\epsilon f(x)\} \\
&= \max\{G \bullet D : G \in \delta_\epsilon \lambda_1(X)\} \\
&= (\tilde{\lambda}_1)'_\epsilon(X; D)
\end{aligned}$$

The advantage of $\delta_\epsilon f(x)$ over the larger $\epsilon$-subdifferential $\partial_\epsilon f(x)$ is that an explicit formula is available (cf. [13], [58, Prop. 3]):

$$\tilde{f}'_\epsilon(x; d) = \sigma_{\delta_\epsilon f(x)}(d) = \lambda_1\left(Q_\epsilon^\top (\mathcal{F}'(x)d)Q_\epsilon\right) = \lambda_1\left(Q_\epsilon^\top D Q_\epsilon\right).$$

One of the main contributions from the work [58] is the following

THEOREM 2.1. *[58, Thm. 4] Let $\epsilon \geq 0$, $\eta \geq 0$ and $X \in \mathbb{S}^m$ and define*

$$(2.3) \qquad \rho(\epsilon, \eta) = \left(\frac{2\eta}{\Delta_\epsilon(X)}\right)^{1/2} + \frac{2\eta}{\Delta_\epsilon(X)}.$$

*Then*

$$(2.4) \qquad \partial_\eta \lambda_1(X) \subset \delta_\epsilon \lambda_1(X) + \rho(\epsilon, \eta) B,$$

*where $B$ is the unit ball in $\mathbb{S}^m$.*

An immediate consequence of this theorem is the estimate:

$$(2.5) \qquad f'_\eta(x; d) \leq \tilde{f}'_\epsilon(x; d) + \rho(\epsilon, \eta)\|D\|.$$

**2.3. Steepest descent.** The direction of steepest descent plays an important role in the analysis of smooth functions. In the case of a non-smooth function it may be obtained by solving the program

$$\min_{\|d\| \le 1} f'(x; d).$$

Fenchel duality shows that

$$\min_{\|d\| \le 1} f'(x; d) = \min_{\|d\| \le 1} \max_{g \in \partial f(x)} g^\top d$$

$$= \max_{g \in \partial f(x)} \min_{\|d\| \le 1} g^\top d$$

$$= \max_{g \in \partial f(x)} -g^\top \frac{g}{\|g\|},$$

so the direction $d$ of steepest descent is obtained as the solution of a convex program:

$$d = -\frac{g}{\|g\|}, \qquad g = \mathrm{argmin}\{\|g\| : g \in \partial f(x)\}.$$

As in the classical case, $d$ will be a direction of descent if there is any, and in that case the relation

$$f'(x; d) = -\|g\| = -\mathrm{dist}(0, \partial f(x)) < 0$$

is satisfied. What is important is that the very same conclusions hold for the $\epsilon$-subdifferential and the $\epsilon$-enlarged subdifferential.

DEFINITION 2.2. *If* $0 \notin \partial_\epsilon f(x)$, *the direction of steepest* $\epsilon$-*descent* $d$ *is*

$$d = -\frac{g}{\|g\|}, \qquad g = \mathrm{argmin}\{\|g\| : g \in \partial_\epsilon f(x)\}$$

*and satisfies* $f'_\epsilon(x; d) = -\mathrm{dist}(0, \partial_\epsilon f(x)) < 0$. *Similarly, if* $0 \notin \delta_\epsilon f(x)$, *then the direction of steepest* $\epsilon$-*enlarged descent is*

$$(2.6) \qquad d = -\frac{g}{\|g\|}, \qquad g = \mathrm{argmin}\{\|g\| : g \in \delta_\epsilon f(x)\}$$

*and satisfies* $\tilde{f}'_\epsilon(x; d) = -\mathrm{dist}(0, \delta_\epsilon f(x)) < 0$.

For a practical implementation it will be convenient to accept approximate solutions of program (2.6). We have the following

LEMMA 2.3. *Let* $0 < \omega \le 1$ *and* $0 \notin \delta_\epsilon f(x)$, *and consider a direction* $d$ *which solves (2.6) approximately in the sense that*

$$(2.7) \qquad d = -\frac{g}{\|g\|}, \qquad \tilde{f}'_\epsilon(x; d) \le -\omega \|g\|.$$

*Then*

$$(2.8) \qquad -\mathrm{dist}(0, \delta_\epsilon f(x)) \le \tilde{f}'_\epsilon(x; d) \le -\omega \, \mathrm{dist}(0, \delta_\epsilon f(x)).$$

*Proof.* Let $\tilde{d} = -\tilde{g}/\|\tilde{g}\|$ be the solution of (2.6), then $\|\tilde{g}\| = \mathrm{dist}(0, \delta_\epsilon f(x))$ and $\tilde{f}'_\epsilon(x; \tilde{d}) = -\|\tilde{g}\|$. As $g \in \delta_\epsilon f(x)$, we have $\|\tilde{g}\| \le \|g\|$. Therefore (2.7) gives $\tilde{f}'_\epsilon(x; d) \le -\omega \|g\| \le -\omega \|\tilde{g}\| = -\omega \, \mathrm{dist}(0, \delta_\epsilon f(x)) < 0$. The left hand estimate follows from $\tilde{f}'_\epsilon(x; d) = \max\{g^\top d : g \in \delta_\epsilon f(x)\} \ge \tilde{g}^\top d \ge -\|\tilde{g}\| = -\mathrm{dist}(0, \delta_\epsilon f(x))$. $\square$

**3. First-order analysis.** In this chapter we derive and analyze a first order bundle algorithm for minimizing non-convex maximum eigenvalue functions $f = \lambda_1 \circ \mathcal{F}$ for $\mathcal{C}^2$-operators $\mathcal{F}$. Occasionally we will specify $\mathcal{F}$ to a bilinear $\mathcal{B}$ or even to an affine $\mathcal{A}$.

**3.1. Optimality conditions.** It is well-known that the $\epsilon$-subdifferential may be used to obtain approximate optimality conditions, which lead to finite termination criteria in a convex minimization algorithm. See [30, 42, 43, 34] on how this is done. How about the meaning of an approximate optimality condition like $0 \in \partial_\epsilon f(x)$ for non-convex maximum eigenvalue function $f$? It is not surprising that without convexity, the consequences of $0 \in \partial_\epsilon f(x)$ are weaker:

LEMMA 3.1. *Let $f = \lambda_1 \circ \mathcal{B}$ with $\mathcal{B}$ bilinear. Let*

$$c := \max_{\|d\|=1} \max_{Z \in C_m} \left| \sum_{i<j} Z \bullet B_{ij} d_i d_j \right|,$$

*where $\mathcal{C}_m = \{ Z \in \mathbb{S}^m : Z \succeq 0, \, \mathrm{tr} Z = 1 \}$. Let $\theta > 0$, $\epsilon \geq 0$, $\sigma \geq 0$ and define $r = r(\epsilon, \sigma, \theta, c)$ as*

$$r = \frac{-\sigma + \sqrt{\sigma^2 + 4\theta c \epsilon}}{2c}.$$

*Then every $x$ such that $\mathrm{dist}(0, \partial_\epsilon f(x)) \leq \sigma$ is $(1+\theta)\epsilon$-optimal within the ball $B(x,r) := \{ x' \in \mathbb{R}^n : \|x' - x\| \leq r \}$. That means, $f(x) \leq \min_{x' \in B(x,r)} f(x') + (1+\theta)\epsilon$.*

*Proof.* Let $X = \mathcal{B}(x)$. By assumption there exists $G \in \partial_\epsilon \lambda_1(X)$ such that $g = \mathcal{B}'(x)^* G$ satisfies $\|g\| \leq \sigma$. Now let $x' \in B(x,r)$ be written as $x' = x + td$ with $\|d\| = 1$, $t = \|x - x'\|$. Put $X' = \mathcal{B}(x')$. By definition of the $\epsilon$-subdifferential, $G \bullet X' \leq \lambda_1(X')$ and $G \bullet X \geq \lambda_1(X) - \epsilon = f(x) - \epsilon$. Therefore

$$
\begin{aligned}
f(x) &\leq G \bullet X + \epsilon \\
&= G \bullet X' + G \bullet (X - X') + \epsilon \\
&\leq \lambda_1(X') + G \bullet (\mathcal{B}(x) - \mathcal{B}(x')) + \epsilon \\
&= f(x') + G \bullet \left( t\mathcal{B}'(x)d + t^2[d, \mathcal{B}''d] \right) + \epsilon \\
&\leq f(x') + \sigma \|x - x'\| + c \|x - x'\|^2 + \epsilon \\
&\leq f(x') + \sigma r + c r^2 + \epsilon = f(x') + (1+\theta)\epsilon
\end{aligned}
$$

by the definition of $r$. This proves the claim. ∎

The result includes several limiting cases. Clearly $\sigma = 0$ is interesting, where we get $r(\epsilon, 0, \theta, c) = \left( \theta \epsilon c^{-1} \right)^{1/2}$. The case $c = 0$ is also possible. It corresponds to an affine operator $\mathcal{A}$, where the second derivative $\mathcal{B}''$ vanishes. We obtain $r(\epsilon, 0, \theta, 0) = +\infty$, meaning that the estimate is a global one. A third limiting case is when $\sigma > 0$ and $c = 0$. Here $r(\epsilon, \sigma, \theta, 0) = \theta \epsilon \sigma^{-1}$. Notice however that in the case $c < \infty$ the look-ahead character of this result is limited. Even while it is true that for $\sigma = 0$ the radius $r \sim \epsilon^{1/2}$ decreases slower than the discrepancy in the values, which behaves like $\sim \epsilon$, we have no way to know whether a local minimum is within the horizon $r$ of our current iterate $x$.

LEMMA 3.2. *Let $f = \lambda_1 \circ \mathcal{F}$. Suppose $\mathrm{dist}(0, \partial_{\epsilon_k} f(x_k)) \to 0$, $x_k \to \bar{x}$ and $\epsilon_k \to 0$. Then $0 \in \partial f(\bar{x})$.*

*Proof.* The proof is straightforward. By the definitions we have $G_k \in \partial_{\epsilon_k} \lambda_1(X_k)$ for some $G_k$ such that $\mathcal{F}'(x_k)^* G_k = g_k \to 0$. In particular, $G_k \bullet X \leq \lambda_1(X)$ for

6

every $X \in \mathbb{S}^m$ and $G_k \bullet X_k \geq \lambda_1(X_k) - \epsilon_k$. Passing to the limit in every convergent subsequence of $G_k$, say $G_k \to \bar{G}$, gives $\bar{G} \bullet X \leq \lambda_1(X)$ for every $X$ and $\bar{G} \bullet \bar{X} = \lambda_1(\bar{X})$, so $\bar{G} \in \partial \lambda_1(\bar{X})$. By the continuity of $\mathcal{F}'$ we have $\mathcal{F}'(\bar{x})^* \bar{G} = 0$. That means $0 \in \partial f(\bar{x})$.
□

Taken together these two Lemmas justify a stopping test based on the smallness of $\mathrm{dist}(0, \partial_\epsilon f(x))$ in tandem with the smallness of $\epsilon$. Such a test is used e.g. in the non-convex bundle method [20], where the Goldstein $\epsilon$-subdifferential is used.

Let us fix $x$ and $\epsilon$ and see what happens when $0 \notin \delta_\epsilon f(x)$. Choose $0 < \omega \leq 1$ and suppose $d$ is an approximate direction of steepest $\epsilon$-enlarged descent satisfying (2.7). The function $\eta \to \rho(\epsilon, \eta)$ in (2.3) is unbounded and monotonically increasing on $[0, \infty)$. Therefore if $\mathcal{F}'(x)d \neq 0$, there exists a unique $\eta = \eta(\epsilon)$ such that $\rho(\epsilon, \eta(\epsilon)) \|\mathcal{F}'(x)d\| = -\frac{1}{2}\tilde{f}'_\epsilon(x; d) > 0$. Using (2.5) gives the following consequence of Theorem 2.1:

$$(3.1) \qquad f'_{\eta(\epsilon)}(x; d) \leq \frac{1}{2}\tilde{f}'_\epsilon(x; d) < 0.$$

The same estimate also holds in the case $\mathcal{F}'(x)d = 0$. Straightforward calculus with (2.3) now shows that in the case $\mathcal{F}'(x)d \neq 0$

$$(3.2) \qquad \eta(\epsilon) = \frac{\Delta_\epsilon(X)}{8}\left(-1 + \sqrt{1 - \frac{2\,\tilde{f}'_\epsilon(x; d)}{\|\mathcal{F}'(x)d\|}}\right)^2.$$

We summarize these observations in the following

LEMMA 3.3. *Suppose $0 \notin \delta_\epsilon f(x)$ for some $\epsilon > 0$. Let $0 < \omega \leq 1$ and let $d$ be a direction of approximate steepest $\epsilon$-enlarged descent satisfying (2.7), (2.8). Then choosing $\eta = \eta(\epsilon)$ as in (3.2) gives $f'_\eta(x; d) \leq \frac{1}{2}\tilde{f}'_\epsilon(x; d) < 0$, i.e., $d$ is a direction of $\eta$-descent. Moreover, $\eta(\epsilon)$ satisfies the estimate*

$$(3.3) \quad \eta(\epsilon) \geq \begin{cases} \dfrac{\Delta_\epsilon(X)\,\omega^2}{18\|\mathcal{F}'(x)d\|^2}\,\mathrm{dist}(0, \delta_\epsilon f(x))^2, & for \;\; 0 > \tilde{f}'_\epsilon(x; d) \geq -3/2\|\mathcal{F}'(x)d\| \\ \dfrac{\Delta_\epsilon(X)\,\omega}{16\|\mathcal{F}'(x)d\|}\,\mathrm{dist}(0, \delta_\epsilon f(x)), & for \;\; \tilde{f}'_\epsilon(x; d) \leq -3/2\|\mathcal{F}'(x)d\| \end{cases}$$

**3.2. Actual and predicted decrease.** Our approach to (1.3) is to estimate the decrease of $f$ in direction $d$ with the help of the convex model $t \mapsto \lambda_1(X + tD)$. The following definition will be helpful.

DEFINITION 3.4. *Let $f = \lambda_1 \circ \mathcal{F}$, and fix $x, d \in \mathbb{R}^n$. Then $\alpha_t = f(x + td) - f(x)$ is called the actual decrease of $f$ at $x$ in direction $d$ with step $t$, while $\pi_t = \lambda_1(X + tD) - f(x)$ is called the predicted decrease at $x$ in direction $d$ with step $t$.*

Naturally, a true decrease of $f$ in direction $d$ with step $t$ only occurs when $\alpha_t < 0$, and similarly for $\pi_t$. During the following, we will use the interplay between $\pi_t$ and $\alpha_t$ in order to find suitable steps $t$ which allow us to quantify $\alpha_t$. We have the formula

$$\begin{aligned} \alpha_t &= f(x + td) - f(x) \\ &= f(x + td) - \lambda_1(X + tD) + \lambda_1(X + tD) - f(x) \\ &= f(x + td) - \lambda_1(X + tD) + \pi_t \end{aligned}$$

so for non-convex $f$ we will have to estimate the mismatch $\alpha_t - \pi_t = f(x + td) - \lambda_1(X + tD)$.

Expanding the $\mathcal{C}^2$ operator $\mathcal{F}$ in a neighborhood of $x$ gives $\mathcal{F}(x + td) = X + tD + t^2H + t^2K_t$, where $K_t \to 0$ as $t \to 0$. In the case of a bilinear $\mathcal{B}$, $K_t = 0$ and $H = [d, \mathcal{B}''d]$ is independent of $x$. Then $f(x+td) = \lambda_1(X+tD+t^2(H+K_t))$ and by Weyl's theorem we obtain

$$t^2\lambda_m(H + K_t) \leq \lambda_1(X + tD + t^2(H + K_t)) - \lambda_1(X + tD) \leq t^2\lambda_1(H + K_t).$$

Altogether, $|f(x + td) - \lambda_1(X + tD)| \leq t^2\|H + K_t\|$. This motivates the following

DEFINITION 3.5. *Let $x, d \in \mathbb{R}^n$, $\|d\| = 1$ and $0 < t < +\infty$ and expand $\mathcal{F}$ as above. Then $L_{x,d,t} := \sup\{\|H + K_\tau\| : 0 \leq \tau \leq t\} < +\infty$. When $\mathcal{F}$ is of class $\mathcal{C}_b^2$, then $t = +\infty$ is allowed and the $L_{x,d,\infty}$ are uniformly bounded on every bounded set of $x$. For bilinear $\mathcal{B}$, $L_{x,d,t} = L_d = \|H\| = \|[d, \mathcal{B}''d]\|$ is independent of $x$ and $t > 0$ and therefore uniformly bounded.*

We summarize our finding by the following

LEMMA 3.6. *Let $0 < T \leq \infty$ such that $L_{x,d,T} < \infty$. Then for every $t \leq T$ the actual decrease $\alpha_t$ satisfies*

$$(3.4) \qquad\qquad \alpha_t = \ell\, t^2 + \pi_t \quad \text{for some } |\ell| \leq L_{x,d,T}.$$

*If $\mathcal{F}$ is of class $\mathcal{C}_b^2$, we may choose $T = \infty$. In the bilinear case $L_{x,d,T} = L_d = \|[d, \mathcal{B}''d]\|$, and in the affine case, $L_{x,d,T} = L_d = 0$.*

Our next step is to quantify the predicted decrease $\pi_t$ of the convex model $t \mapsto \lambda_1(X + tD)$. This is done in the next section.

**3.3. Directional analysis.** Let $x$ be the current iterate in a tentative algorithm, and let $d$ with $\|d\| = 1$ be a direction such that for some $\eta > 0$, $f'_\eta(x; d) < 0$, i.e., $d$ is a direction of $\eta$-descent at $x$. This could happen with $\eta = \eta(\epsilon)$ and $d$ an approximate direction of steepest $\epsilon$-enlarged descent as in (2.7). Then $0 \notin \partial_\eta f(x)$. As usual let $X = \mathcal{F}(x)$ and $D = \mathcal{F}'(x)d$, then by definition, $f'_\eta(x; d) = (\lambda_1)'_\eta(X; D)$, so $(\lambda_1)'_\eta(X; D) < 0$. Following [30, XI.1], there exists a hyperplane supporting the epigraph of $\lambda_1$, which passes through the point $(X, \lambda_1(X) - \eta)$ and touches the epigraph of $\lambda_1$ at a point $(X + t_\eta D, \lambda_1(X + t_\eta D))$, except when $t \mapsto \lambda_1(X + tD)$ is affine on $[0, \infty)$, or has an asymptote with slope $f'_\eta(x; d)$. In those cases let $t_\eta = \infty$ for consistency. If there are several steps with this property, then for definiteness let $t_\eta$ be the smallest one. We shall say that the step $t_\eta$ *realizes* the $\eta$-directional derivative. Notice that $\lambda'_1(X + t_\eta D; D) \leq (\lambda_1)'_\eta(X; D) \leq -\lambda'_1(X + t_\eta D; -D)$, so for almost all $t_\eta$ we have equality $(\lambda_1)'_\epsilon(X; D) = \lambda'_1(X + t_\eta D; D)$. In general there exists at least a subgradient $G \in \partial\lambda_1(X + t_\eta D)$ such that $(\lambda_1)'_\eta(X; D) = G \bullet D$.

All this being a purely directional situation, we could also describe the case by introducing the function $\phi(t) = \lambda_1(X + tD)$. Then the line passing through $(0, \phi(0) - \eta)$ touches the epigraph of $\phi$ at $(t_\eta, \phi(t_\eta))$. The reader may want to inspect Figure 2.1.2 on page 105 of [30] for some illustration of these on-goings.

For the following assume that $\phi$ is not affine on $[0, \infty)$. The case $t_\eta = \infty$ with an asymptote is allowed. Suppose we backtrack and consider all lines passing through $(0, \phi(0) - \eta')$ for some $0 \leq \eta' \leq \eta$, touching the epigraph of $\phi$ at the corresponding points $(t_{\eta'}, \phi(t_{\eta'}))$, where again $t_{\eta'}$ is the smallest step if there are several. Then we introduce a function $\eta' \to t_{\eta'}$ with the following properties: $t_0 = 0$, and $t_\eta = t_{\eta'}$ at $\eta' = \eta$. It is monotonically increasing by convexity of $\phi$. (Notice that its inverse, $t \to \eta(t)$, is not a function in the strict sense unless $\phi$ is differentiable. But we can introduce the notation to indicate any choice such that $\eta(t(\eta)) = \eta$.)

Recall that we assume $f'_\eta(x; d) < 0$. Then the decrease of $\phi$ on $[0, t_\eta]$ is

$$(3.5) \qquad\qquad \phi(t_\eta) - \phi(0) = -\eta + t_\eta f'_\eta(x; d) < -\eta < 0.$$

8

Therefore the secant joining the points $(0, \phi(0))$ and $(t_\eta, \phi(t_\eta))$ has slope $-\sigma$, where

$$(3.6) \qquad \sigma = \frac{\eta - t_\eta f'_\eta(x; d)}{t_\eta} = \frac{\eta}{t_\eta} - f'_\eta(x; d) > 0.$$

We can then get a pessimistic estimate of $\pi_t$ by using the weaker decrease of the secant, which for a step $t \leq t_\eta$ decreases by $-\sigma t$. Therefore $\alpha_t = \ell t^2 + \pi_t \leq L t^2 - \sigma t$ for every $t \leq t_\eta$, where $L := L_{x,d,t_\eta}$. (If $t_\eta = \infty$ and $L_{x,d,\infty} = \infty$, then this estimate will only be true for some finite $T$, $L := L_{x,d,T}$ and all $t \leq T$.) The minimum of the term $L t^2 - \sigma t$ on $[0, t_\eta]$ is attained either at $t = \sigma/2L$ if $\sigma/2L \leq t_\eta$, or at $t = t_\eta$. Substituting into $\alpha_t$ gives the following

LEMMA 3.7. *Suppose $t_\eta < \infty$. Let $d$ be an approximate direction of steepest $\epsilon$-enlarged descent satisfying (2.7) and let $\eta = \eta(\epsilon)$ be as in (3.2), so $f'_\eta(x; d) < 0$.*
*(a) If $\sigma$ defined by (3.6) satisfies $\sigma \leq 2 L_{x,d,t_\eta} t_\eta$, then the step $t_\sigma := \sigma/2 L_{x,d,t_\eta}$ gives a guaranteed decrease*

$$(3.7) \qquad \alpha_{t_\sigma} \leq -\frac{\sigma^2}{4 L_{x,d,t_\eta}} = -\frac{1}{4 L_{x,d,t_\eta}} \left( \frac{\eta}{t_\eta} - f'_\eta(x; d) \right)^2$$

$$\leq -\frac{1}{4 L_{x,d,t_\eta}} \left( \frac{\eta}{t_\eta} - \frac{1}{2} \tilde{f}'_\epsilon(x; d) \right)^2 < 0.$$

*(b) If on the other hand $\sigma > 2 L_{x,d,t_\eta} t_\eta$, then the step $t_\eta$ guarantees the decrease*

$$(3.8) \qquad \alpha_{t_\eta} \leq -\frac{\sigma t_\eta}{2} = \frac{1}{2} \pi_{t_\eta} = -\frac{1}{2} \left( \eta - t_\eta f'_\eta(x; d) \right) < 0.$$

At first sight we would probably prefer case (b) over case (a), as it seems to give a better rate of decrease $\mathcal{O}(\eta)$ versus $\mathcal{O}(\eta^2)$. Moreover, the constant $L$ is likely to be large, so $\alpha_{t_\sigma}$ is expected to be small. In the same vein, $t_\eta$ could be large, so the main contribution in $\alpha_{t_\sigma}$ probably comes from the term $f'_\eta(x; d)$ respectively $\tilde{f}'_{\epsilon_k}(x_k; d_k)$, which also occurs with order 2, as opposed to the second branch in (3.2). But (a) has a surprising advantage over (b). Namely, when later on our algorithm will take steps $t_k$ which force the terms $\alpha_{t_k}$ to tend to zero, $t_k = t_\sigma$ in case (a) will imply $\tilde{f}'_{\epsilon_k}(x_k; d_k) \to 0$. On the other hand, $t_k = t_\eta$ in case (b) will only lead to $\Delta_{\epsilon_k}(X_k) \tilde{f}'_{\epsilon_k}(x_k; d_k) \to 0$. Dealing with the term $\Delta_{\epsilon_k}(X_k)$ will cause some trouble. Notice, however, that as a rule we have to expect case (b). In particular, for convex $f = \lambda_1 \circ \mathcal{A}$ we have $L = 0$, so case (b) is always on (compare the discussion in [58]).

Let us catch up with the case where $t_\eta = \infty$. This may happen if $\phi : t \mapsto \lambda_1(X + tD)$ is affine on $[0, \infty)$ or has an asymptote with slope $f'_\eta(x; d) < 0$. Then $\phi$ is below the line with slope $-\sigma = f'_\eta(x; d)$ passing through the point $(0, \phi(0))$, and we may use this line to estimate $\pi_t$. The result is the same as in Lemma 3.7, with $t_\eta = \infty$ and $\eta/t_\eta = 0$.

LEMMA 3.8. *Suppose $f'_\eta(x; d) < 0$ and $t_\eta = \infty$. Fix $T > 0$ such that $L := L_{x,d,T} < \infty$. Then the largest possible decrease of $f$ in direction $d$ amongst steps of length $t \leq T$ is obtained at $t_\sigma = \sigma/2L$ with $\alpha_{t_\sigma} \leq -\sigma^2/4L$ if $t_\sigma \leq T$, otherwise at $t = T$ with $\alpha_t \leq LT^2 - \sigma T \leq -\frac{1}{2} T \sigma$.*

For $\mathcal{F}$ of class $\mathcal{C}_b^2$ we let $T = \infty$. In particular, this works for bilinear $\mathcal{B}$, so there is no restriction on the steplength $t_\sigma$. For general $\mathcal{C}^2$ operators the result seems to pose a little problem, as we need to know $T$ to compute $L = L_{x,d,T}$, and $L$ in order

9

to check whether $t_\sigma \leq T$. Notice, however, that $L_{x,d,T}T$ increases as $T \to \infty$, while $t_\sigma > T$ for all $T$ meant $L_{x,d,T}T$ remained bounded by $\sigma/2$. This is only possible when $L_{x,d,T} = 0$. So there is always the possibility to escape this dilemma. But section 3.8 will show an even simpler way to deal with general $\mathcal{C}^2$ operators.

**3.4. Decrease of the order $\mathcal{O}(\eta)$.** Let us look more systematically at those cases where a quantifiable decrease of the order $\mathcal{O}(\eta)$ is possible. This requires a sufficiently small trial step $t_\eta$. We have the following

LEMMA 3.9. *Let $0 < \rho_0 < 1$. Then the trial step $t_\eta$ gives a decrease $\alpha_{t_\eta} \leq -(1 - \rho_0)\eta$ provided $t_\eta \leq \vartheta$, where $\vartheta$ is the critical stepsize*

$$(3.9) \qquad \vartheta := \frac{-f'_\eta(x;d) + \sqrt{f'_\eta(x;d)^2 + 4L\rho_0\eta}}{2L}.$$

*Here $L = L_{x,d,t_\eta}$, and the limiting case $L_{x,d,t_\eta} = 0$ is allowed and gives $\vartheta = \infty$.*

*Proof.* Since $|\ell| \leq L$, the inequality $\alpha_{t_\eta} = \ell t_\eta^2 - \eta + t_\eta f'_\eta(x;d) \leq -(1 - \rho_0)\eta$ is satisfied as soon as the stronger quadratic inequality $Lt_\eta^2 + t_\eta f'_\eta(x;d) - \rho_0\eta \leq 0$ holds. The corresponding quadratic equation has two real solutions, and the positive solution is $\vartheta$ in (3.9). Therefore the quadratic inequality holds as soon as $t_\eta \leq \vartheta$. □

**3.5. The $\epsilon$-management.** The dependence of the estimate (3.2) on the gap $\Delta_\epsilon(X)$ suggests that we choose $\Delta_\epsilon(X)$ as large as possible. This strategy is indeed best when we aim at a finite termination theorem (see [13], [58]). If we want to prove convergence, the situation is, as we shall see, more subtle. Here we want $\Delta_\epsilon(X)$ large, but with the proviso that $\epsilon \to 0$.

LEMMA 3.10. *Let $\bar{\epsilon} > 0$. Then there exists $\epsilon \leq \bar{\epsilon}$ such that $\Delta_\epsilon(X) \geq \bar{\epsilon}/m$.*

*Proof.* By definition of $r(\bar{\epsilon})$ we have $\lambda_1(X) \geq \cdots \geq \lambda_{r(\bar{\epsilon})}(X) \geq \lambda_1(X) - \bar{\epsilon} > \lambda_{r(\bar{\epsilon})+1}(X)$. That means the $r(\bar{\epsilon})$ gaps $\lambda_i(X) - \lambda_{i+1}(X)$, $i = 1, \ldots, r(\bar{\epsilon})$, add up to $\lambda_1(X) - \lambda_{r(\bar{\epsilon})+1}(X)$, which exceeds $\bar{\epsilon}$. So at least one of these gaps is larger than $\bar{\epsilon}/r(\bar{\epsilon})$, hence larger than $\bar{\epsilon}/m$. Suppose a gap exceeding $\bar{\epsilon}/m$ is $\lambda_i(X) - \lambda_{i+1}(X)$. Then we put $\epsilon = \lambda_1(X) - \lambda_i(X)$. □

In practice we could either choose $i$ smallest possible with $\lambda_i(X) - \lambda_{i+1}(X) \geq \bar{\epsilon}/m$, or we could pick $i$ so that $\lambda_i(X) - \lambda_{i+1}(X)$ is largest possible. In the first case we get a small $\epsilon$, which reduces the numerical burden to compute $Q_\epsilon$. In the second case we render (3.2) the most efficient, possibly with a larger $\epsilon$.

DEFINITION 3.11. *The maximum eigenvalue function $f = \lambda_1 \circ \mathcal{F}$ is called linearly bounded below if for every $x \in \mathbb{R}^n$, the mapping $e \mapsto \lambda_1(\mathcal{F}(x) + \mathcal{F}'(x)e)$ is bounded below.*

The significance of this definition is in the following

LEMMA 3.12. *Suppose $f = \lambda_1 \circ \mathcal{F}$ is linearly bounded below. Let $x \in \mathbb{R}^n$ be such that $0 \notin \partial f(x)$. Then there exists $\bar{\epsilon} > 0$ such that $0 \in \delta_{\bar{\epsilon}} f(x)$ but $0 \notin \delta_\epsilon f(x)$ for all $0 \leq \epsilon < \bar{\epsilon}$.*

*Proof.* Let $\bar{\epsilon} = \lambda_1(X) - \lambda_m(X)$, then $\delta_{\bar{\epsilon}} f(x) = \mathcal{F}'(x)^* (\mathcal{C}_m)$. By [58, Lemma 6] we have $0 \in \mathcal{F}'(x)^* (\mathcal{C}_m)$ because $e \mapsto \lambda_1(\mathcal{F}(x) + \mathcal{F}'(x)e)$ is bounded below. So we have proved $0 \in \delta_{\bar{\epsilon}} f(x)$. On the other hand, $0 \notin \delta_\epsilon f(x)$ if $\epsilon$ is sufficiently small. Take for instance $\epsilon$ so small that $r(\epsilon)$ equals the multiplicity of $\lambda_1(X)$. Then $\delta_\epsilon f(x) = \partial f(x)$, hence $0 \notin \delta_\epsilon f(x)$ by our hypothesis. By reducing $\bar{\epsilon}$ we can make it smallest possible with $0 \in \delta_{\bar{\epsilon}} f(x)$. □

It is easy to force a maximum eigenvalue function to be linearly bounded below on a bounded set $C$. More generally, if $f$ is bounded below by $\gamma \in \mathbb{R}$ on a set $C$, then define $\tilde{\mathcal{F}}(x) = \text{diag}(\gamma, \mathcal{F}(x)) \in \mathbb{S}^{m+1}$. Clearly $\tilde{f} = \lambda_1 \circ \tilde{\mathcal{F}}$ agrees with $f$ on $C$

10

and is linearly bounded below. During the following we will assume that $f$ is linearly bounded below on the level set $\{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$.

**3.6. Line search.** We need one more element before our first-order bundle algorithm may be presented in detail. The analysis so far suggests several steps $t$ which give a quantifiable decrease $\alpha_t$. But we have to make sure that such a step is found by a finite procedure.

Suppose the parameter $\epsilon$ has been chosen, $\tilde{f}'_\epsilon(x; d) < 0$ has been computed and $\eta = \eta(\epsilon)$ has been found as in (3.2). Then we need to find $t_\eta$ realizing $f'_\eta(x; d)$. If done in a precise way, this search is likely to be very costly. We therefore choose a relaxation.

Fix a tolerance parameter $0 < \theta_0 < 1$ and seek $t > 0$ such that $\lambda'_1(X + tD; D) < \frac{1}{4}\tilde{f}'_\epsilon(x; d) < 0$ and $\eta(t) > \theta_0 \eta$. The set of these $t$ is nonempty and contains $t_\eta$ as an interior point. This means a line search procedure like bisection can find $t > 0$ with these properties in a finite number of steps. So altogether replacing the original $\eta$ by $\theta_0 \eta$ still gives the same order of decrease, but has the benefit to locate an approximation of the realizing abscissa in a finite procedure. This relaxation means that some of the estimates in previous Lemmas will get an extra factor $\theta_0$ or $\theta_0^2$.

Naturally, while trying to locate $t_\eta$, we should not forget our original purpose of reducing $f$. After all, $t_\eta$ relates only to $t \mapsto \lambda_1(X + tD)$, and not directly to $f$. We should therefore evaluate $\alpha_t$ at each intermediate step $t$ visited during the search for $t_\eta$. Using the parameter $\rho_0 \in (0, 1)$ from section 3.4, we may accept an intermediate $t$ immediately if $\alpha_t \leq -(1 - \rho_0)\eta$, because this is the best order of decrease we can hope to achieve in general.

Suppose $t_\eta$ respectively its substitute has been found. Keep among the intermediate steps the one with maximum decrease $\alpha_t$ and call it $\zeta$. If necessary continue and compute the decrease $\alpha_{t_\sigma}$ at $t_\sigma$, where $\sigma$ is as in (3.6) and compare $\alpha_{t_\sigma}$ to $\alpha_\zeta$. This covers all possible cases. For instance, if $t_\eta < \vartheta$ with $\vartheta$ as in (3.9), we will get a decrease of the order $\mathcal{O}(\eta)$.

Except in the bilinear case, we need to estimate $L_{x,d,t}$ for various stepsizes $t$. Strictly speaking this could not be done in a finite procedure. However, $t$ and therefore $\|K_t\|$ are expected to be small, so $L_{x,d,t} \sim \|H\|$ for sufficiently small $t$, meaning that good estimates are available. In any case, the constants $L_{x,d,t}$ need not be known exactly. An upper bound $L$ for all the $L_{x,d,t}$ with $x$ ranging over the level set $\{x : f(x) \leq f(x_0)\}$ is all that is needed to prove convergence.

**3.7. First order algorithm.** In this section we present our first order bundle algorithm for the unconstrained minimization of $f = \lambda_1 \circ \mathcal{F}$ (see Figure 1, next page) and prove convergence.

Let us consider sequences $x_k$, $d_k$ and $\epsilon_k^\sharp, \bar{\epsilon}_k, \epsilon_k$ generated by the bundle algorithm. Suppose the sequence $x_k$ is bounded, the $f(x_k)$ are bounded below, and that $\mathcal{F}$ is of class $\mathcal{C}_b^2$. Then $f(x_j) - f(x_0) = \sum_{k=0}^{j-1} f(x_{k+1}) - f(x_k) = \sum_{k=0}^{j-1} \alpha_{t_k}$ is bounded, and since each coefficient $\alpha_{t_k}$ is negative, the series $\sum_{k=0}^{\infty} \alpha_{t_k}$ converges. By boundedness of $\mathcal{F}''$ and the $x_k$, the constants $L_{x_k,d_k,t_{\eta_k}}$ are uniformly bounded. Then Lemma 3.7 and the choice of $t_k$ in step 8 give

$$\alpha_{t_k} \leq -K \max\left\{ \left(\eta_k/t_{\eta_k} - \tilde{f}'_{\epsilon_k}(x_k; d_k)\right)^2, \eta_k \right\} < 0$$

for some $K > 0$ independent of $k$. This implies either $\tilde{f}'_{\epsilon_k}(x_k; d_k) \to 0$ or $\Delta_{\epsilon_k}(X_k)\, \tilde{f}'_{\epsilon_k}(x_k; d_k) \to 0$, depending on which of the cases in Lemma 3.7 occurs. Therefore, by (2.8), and

since $\delta_{\epsilon_k} f(x_k) \subset \partial_{\epsilon_k} f(x_k)$, this implies $\mathrm{dist}(0, \partial_{\epsilon_k} f(x_k)) \to 0$ in the first case, and $\Delta_{\epsilon_k}(X_k) \, \mathrm{dist}(0, \partial_{\epsilon_k} f(x_k)) \to 0$ in the second.

Spectral bundle Algorithm for (1.3)

1. Choose an initial iterate $x_0$ and fix $0 < \theta_0, \rho_0 < 1$ and $0 < \omega \le 1$. Let $\gamma_k > 0$ be a sequence converging slowly to 0. Fix $\epsilon_0^\sharp = \rho_0$. Initialize $S_0 = \emptyset, F_0 = \emptyset$ and let $\texttt{slope} \in \{\texttt{steep}, \texttt{flat}\}$ be a binary variable.

2. Given the current iterate $x_k$, stop if $0 \in \partial f(x_k)$. Otherwise let $\bar{\epsilon}_k > 0$ such that $0 \in \delta_{\bar{\epsilon}_k} f(x_k)$ but such that $0 \notin \delta_\epsilon f(x_k)$ for $\epsilon < \bar{\epsilon}_k$. Choose $\epsilon_k \le \min\{\bar{\epsilon}_k, \epsilon_k^\sharp\}$ such that $\Delta_{\epsilon_k}(X_k) \ge \min\{\bar{\epsilon}_k/m, \epsilon_k^\sharp/m\}$.

3. compute a direction $d_k$ of approximate steepest $\epsilon_k$-enlarged descent
$$d_k = -\frac{g_k}{\|g_k\|}, \qquad \tilde{f}'_{\epsilon_k}(x_k; d_k) \le -\omega \|g_k\|.$$

4. If $|\tilde{f}'_{\epsilon_k}(x_k; d_k)| \le \epsilon_k^\sharp$ put $\texttt{slope = flat}$, otherwise put $\texttt{slope = steep}$.

5. compute $\eta_k = \eta(\epsilon_k)$ according to (3.2).

6. Search for $t_{\eta_k}$ using a backtracking line search. If during the search $t$ satisfying $\alpha_t \le -(1 - \rho_0)\eta_k$ is found, put $t_k = t$ and goto 9. Otherwise stop as soon as $t$ satisfying $\lambda'_1(X_k + tD_k; D_k) < \frac{1}{4}\tilde{f}'_{\epsilon_k}(x_k; d_k)$ and $\eta(t) \ge \theta_0 \eta_k$ is found. Replace $\eta_k$ by $\eta(t)$ and $t_{\eta_k}$ by $t$. Let $\zeta_k$ be the step which gave the best $\alpha_{\zeta_k}$ during the search.

7. compute $L = L_{x_k, d_k, t_{\eta_k}}$, the trial step $t_{\sigma_k}$ in (3.7) and $\alpha_{t_{\sigma_k}}$.

8. compute $\vartheta_k$ by (3.9) and $\alpha_{\vartheta_k}$. Let $t_k \in \{\zeta_k, \vartheta_k, t_{\eta_k}\}$ the step which gives the best decrease.

9. Put $x_{k+1} = x_k + t_k d_k$. If $\texttt{slope == steep}$ let $F_{k+1} = F_k$, $\epsilon_{k+1}^\sharp = \epsilon_k^\sharp$ $S_{k+1} = S_k \cup \{k\}$. If $\texttt{slope == flat}$, let $F_{k+1} = F_k \cup \{k\}$, $S_{k+1} = S_k$ and put $\epsilon_{k+1}^\sharp = \gamma_\ell$, where $\ell = \text{card}(F_{k+1})$. Replace $k$ by $k + 1$ and go back to step 2.

FIGURE 1.

13

Let $F = \bigcup_k F_k$ and $S = \bigcup_k S_k$ be the flat respectively steep iterates. There are two cases. Suppose first that there is an infinite number of flat iterates $k \in F$. Then $\epsilon_k^\sharp \to 0$ by step 9, because $\epsilon_k^\sharp = \gamma_{\ell_k}$, with $\ell_k$ the number of flat steps that occurred among the first $k+1$ steps, so $\ell_k \to \infty$. Therefore also $\epsilon_k \to 0$. On the other hand, the flat steps $k \in F$ satisfy $|\tilde{f}_{\epsilon_k}'(x_k; d_k)| \leq \epsilon_k^\sharp \to 0$. Then $\mathrm{dist}(0, \partial_{\epsilon_k} f(x_k)) \to 0$ because of (2.8). Therefore if $\bar{x}$ is an accumulation point of the flat subsequence $x_k$, $k \in F$, we conclude with Lemma 3.2 that $0 \in \partial f(\bar{x})$.

Let us next assume that all but a finite number of steps are steep, i.e., $k \in S$ for all $k \geq k_0$. By step 4, $|\tilde{f}_{\epsilon_k}'(x_k; d_k)| > \epsilon_k^\sharp$, $k \geq k_0$, and by step 9 the algorithm stops driving $\epsilon_k^\sharp$ to 0. That means $\tilde{f}_{\epsilon_k}'(x_k; d_k)$ stays away from 0. Then we must have $\Delta_{\epsilon_k}(X_k) \to 0$. (In particular, this rules out case (a) in Lemma 3.7.) Now $\min\{\bar{\epsilon}_k, \epsilon_k^\sharp\} = \bar{\epsilon}_k$ eventually, because by step 2 this minimum tends to 0, while $\epsilon_k^\sharp$ stays away from 0. Hence $\bar{\epsilon}_k \to 0$, and since $0 \in \delta_{\bar{\epsilon}_k} f(x_k) \subset \partial_{\bar{\epsilon}_k} f(x_k)$, Lemma 3.2 implies $0 \in \partial f(\bar{x})$ for every accumulation point of the entire sequence of iterates $x_k$. Altogether we have proved the following

THEOREM 3.13. *Let $f = \lambda_1 \circ \mathcal{F}$ be a (non-convex) maximum eigenvalue function with $\mathcal{F}$ of class $\mathcal{C}_b^2$. Let $x_0$ be fixed so that the level set $\{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$ is compact. Suppose the sequence $x_k$, starting with $x_0$, is generated by the first-order bundle algorithm. Then $x_k$ is bounded and the $f(x_k)$ decrease monotonically. If all but finitely many iterates are steep, ($k \in S$, $k \geq k_0$,) then every accumulation point of $x_k$ is a critical point. If $F$ is infinite, then every accumulation point of the flat subsequence $x_k$, $k \in F$, is a critical point of $f$.*

In [58] the author considers the convex case $f = \lambda_1 \circ \mathcal{A}$ and aims at finite termination. The following is therefore a complement to [13] and [58, Thm. 7]:

COROLLARY 3.14. *Suppose $f = \lambda_1 \circ \mathcal{A}$ is convex. Let $x_0$ be such that the level set $\{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$ is compact. Then every accumulation point of the sequence of iterates $x_k$ generated by the first-order bundle algorithm starting with $x_0$ is a minimum of $f$.*

*Proof.* By convexity every critical point $\bar{x}$ of $f$ is a (global) minimum. Therefore, if all but finitely many $k$ are steep, the Theorem tells us that we are done. Suppose then that the subsequence $k \in F$ is infinite. Then by the above $x_k$, $k \in F$, has an accumulation point, $\bar{x}$, which is a minimum of $f$. Since the algorithm is of decent type, potential other accumulation points $\tilde{x}$ of the steep subsequence $k \in S$ satisfy $f(\tilde{x}) = f(\bar{x})$, hence by convexity are also minima. That proves the claim. $\square$

Clearly this argument also applies when $f$ is non-convex and the accumulation point $\bar{x}$ above is a global minimum of $f$.

REMARK. In the convex case the bundle algorithm essentially agrees with that of Cullum et al. [13] and Oustry [58], except that we force $\epsilon_k \to 0$ in certain cases in order to assure convergence. In [27] Helmberg and Rendl propose an alternative way to maintain an approximation $\widehat{\mathcal{W}}$ of $\partial_\epsilon \lambda_1(X)$, using an orthogonal matrix $P$ other than $Q_\epsilon$. Moreover, they include the possibility to remember previous steps via an aggregate subgradient. Both approaches are compared in [26], and some merits of $\widehat{\mathcal{W}}$ are observed. In [27] the authors prove convergence of their method and in 3.1 claim that convergence based on $\delta_\epsilon f(x)$ requires partial knowledge of the multiplicity $\bar{r}$ of $\lambda_1(\bar{X})$ at the limit $\bar{X}$. More precisely, they claim that $r(\epsilon_k) \geq \bar{r}$ at iterates $X_k$ near $\bar{X}$ is needed. While this is true if the $\epsilon$-management from [13, 58] is used, Corollary 3.14 shows that our way of choosing $\epsilon_k$ gives at least subsequence convergence without guessing $\bar{r}$ correctly. (Naturally, if $f = \lambda_1 \circ \mathcal{A}$ has a strict minimum, our method gives

convergence. Helmberg and Rendls' method converges without this hypothesis.)

REMARK. The reader will have understood already that we intend the covering sequence $\gamma_k$ to converge so slowly that the flat case almost never occurs. Nonetheless, it may seem a little puzzling that when $F$ and $S$ are both infinite, it is the flat subsequence which gets the merit of (subsequence) convergence, while the steep subsequence, seemingly doing all the work on the way, does not get rewarded by subsequence convergence. Indeed, the estimates suggest that steps where $|\tilde{f}'_{\epsilon_k}(x_k; d_k)|$ is large give the best decrease in the cost, and those are the ones in $S$.

Let us therefore examine the case where both $F$ and $S$ are infinite more closely. As we shall see, in most cases the subsequence $S$ will still have many convergent subsequences. Notice first that it is possible that $\tilde{f}'_{\epsilon_k}(x_k; d_k) \to 0$, for a subsequence $k \in S' \subset S$, even though the speed is necessarily slower than that of $\epsilon_k^{\sharp} \to 0$. Since $\epsilon_k \to 0$, a consequence of the fact that $F$ is infinite, we conclude in that case (via Lemma 3.2) that every accumulation point $\tilde{x}$ of $S'$ is critical. Altogether, subsequences like $S'$ are welcome.

Let us next examine a subsequence $S'' \subset S$ where $\tilde{f}'_{\epsilon_k}(x_k; d_k) \leq \tau < 0$, $k \in S''$. In that case we know that $\Delta_{\epsilon_k}(X_k) \to 0$, $k \in S''$. In particular, this may not happen if case (a) in Lemma 3.7 is on. Assuming that we are in case (b) of that Lemma, suppose we have $\bar{\epsilon}_k \to 0$, $k \in S''$. Then we are again done, ending up with another good subsequence $S''$ exhibiting subsequence convergence.

So finally the bad case is when $\bar{\epsilon}_k$, $k \in S''$, stay away from 0. Then by step 2 of the algorithm, we will eventually have $\epsilon_k^{\sharp} < \bar{\epsilon}_k$, $k \in S''$. Now we have to remember that $\sum_k \alpha_{t_k}$ is even summable, a fact we have never exploited so far. From estimate (3.2) we deduce that $\sum_{k \in S''} \Delta_{\epsilon_k}(X_k)^2 < \infty$, hence $\sum_{k \in S''}(\epsilon_k^{\sharp})^2 < \infty$ by the above. Put differently, a subsequence $S''$ of this last type must be extremely sparse, because as we agreed, $\gamma_k$ tends to 0 very slowly. We may for instance decide that it converges so slowly that $\sum_k \gamma_k^2 = \infty$. Then also $\sum_k(\epsilon_k^{\sharp})^2 = \infty$.

This observation at least partially resolves the following dilemma caused by our algorithm. Suppose our method proposes iterates $x_k$ with $k \in F$ and $k \in S$ fairly mixed. Then in order to be on the safe side, we would probably stop the process when $k \in F$. But what to do when all the iterates are in $S$? In practice this will be satisfactory, as we are probably in the case where $F$ is finite. But of course we can never be sure, having to stop after a finite number of steps. It will then be reassuring to know that the probability to be in a subsequence $S''$ of the last type, where subsequence convergence may fail, is in some sense very low. For instance, the probability to meet an element of $S''$ in an interval of fixed length $\ell$, say in $I_n = [n, n + \ell]$, will tend to 0, which makes it very likely that stopping the algorithm in $I_n$ gives an iterate belonging to some of the "good" subsequences of $S$.

**3.8. Bounding $t_\eta$.** For general $\mathcal{C}^2$ operators $\mathcal{F}$ our algorithm has to be mildly modified. Here we encounter the problem that the constants $L_{x,d,t_\eta}$ may become arbitrarily large due to the fact that the trial steps $t_\eta$ may become arbitrarily large. In particular, we then cannot allow cases where $t_\eta = \infty$. As we have seen, this does not occur for bilinear $\mathcal{B}$, so we might regard this case as of minor importance for the applications we have in mind. But eigenvalue programs with general $\mathcal{C}^2$ operators $\mathcal{F}$ have frequently been treated in the literature; for applications in automatic control see for instance [3], [4]. We therefore include a discussion of this case here. The way it is handled is by brute force. We oblige the steps $t_\eta$ to be uniformly bounded by a constant $T > 0$.

In order to understand the difficulty, consider the known case of an affine operator $\mathcal{A}$. If the minimization of $f = \lambda_1 \circ \mathcal{A}$ is to be well-defined, i.e., if $f$ is to be bounded below, then the linear part $A$ of $\mathcal{A}$ needs to satisfy $\lambda_1(Ad) \geq 0$ for every $d$. That is, $Ad$ is not negative definite for any $d$. Moreover, if the set of minimizers of $f = \lambda_1 \circ \mathcal{A}$ is to be bounded, we even require coercivity of $f$, which means $\lambda_1(Ad) > 0$ for every $d \neq 0$.

What happens for general non-convex $\mathcal{C}^2$ operators? Here we notice a difference with the affine case. Even when $f = \lambda_1 \circ \mathcal{F}$ is nicely bounded below and coercive, that is $f(x) \to +\infty$ as $\|x\| \to \infty$, there is no reason why its linearizations about a point $x$, that is, $e \mapsto \lambda_1(\mathcal{F}(x) + \mathcal{F}'(x)e)$, should share this property. In fact some of the linearizations may fail to be coercive, and this may lead to arbitrarily large values $t_\eta$ as above. This effect motivates the following

DEFINITION 3.15. *A representation $f = \lambda_1 \circ \mathcal{F}$ of $f$ is called linearly coercive if $e \mapsto \lambda_1(\mathcal{F}(x) + \mathcal{F}'(x)e)$ is coercive for every $x$.*

The meaning of this definition becomes clear with the following

LEMMA 3.16. *Suppose $f = \lambda_1 \circ \mathcal{F}$ is linearly coercive. Then for every $R > 0$ there exists $T = T(R)$ such that for every $\|x\| \leq R$, every direction $\|d\| = 1$ and every $\eta > 0$ having $f'_\eta(x; d) < 0$, the abscissa $t_\eta$ realizing $f'_\eta(x; d)$ satisfies $t_\eta \leq T$.*

*Proof.* Suppose on the contrary that there exist $\|x_k\| \leq R$, $\|d_k\| = 1$ and $\eta_k > 0$ such that $-\lambda_1'(X_k + t_{\eta_k}D_k; -D_k) \leq f'_{\eta_k}(x_k; d_k) < 0$ is satisfied, where $t_{\eta_k}$ realize $f'_{\eta_k}(x_k; d_k)$ and $t_{\eta_k} \to \infty$. Then each of the functions $\phi_k : t \mapsto \lambda_1(X_k + tD_k)$ decreases on the interval $[0, t_{\eta_k}]$. In particular, for every intermediate $0 < t < t_{\eta_k}$ we have $\lambda_1'(X_k + tD_k; D_k) < 0$ and $\lambda_1(X_k) > \lambda_1(X_k + tD_k)$.

Passing to subsequences, we may assume $x_k \to x$, $d_k \to d$, hence $X_k \to X$, $D_k \to D$. Now consider an arbitrary $t > 0$. Then $t_{\eta_k} > t$ for $k$ large enough, hence $\lambda_1(X_k) > \lambda_1(X_k + tD_k)$ for $k$ large enough. Then in the limit, $\lambda_1(X) \geq \lambda_1(X + tD)$. Since $t > 0$ was arbitrary, $t \mapsto \lambda_1(X + tD)$ is bounded above by $\lambda_1(X)$ on $[0, \infty)$. That contradicts linear coercivity. $\square$

The proof is not constructive, so it is not clear at the moment how $T$ should be computed. However, as we shall see, $T$ will not be required explicitly in the algorithm and rather serves as a theoretical parameter to obtain convergence. Let us now see how a given maximum eigenvalue function $f = \lambda_1 \circ \mathcal{F}$ could be forced to linear coercivity.

LEMMA 3.17. *Let $f = \lambda_1 \circ \mathcal{F}$ be given. For $R > 0$ and $\epsilon_0 > 0$, there exists a linearly coercive maximum eigenvalue function $\tilde{f} = \lambda_1 \circ \tilde{\mathcal{F}}$ such that*

*(i) $f(x) = \tilde{f}(x)$ for every $\|x\| \leq R$.*

*(ii) For every $0 \leq \epsilon \leq \epsilon_0$ and $\|x\| \leq R$, the $\epsilon$-enlarged subdifferentials $\delta_\epsilon f(x)$ and $\delta_\epsilon \tilde{f}(x)$ coincide.*

*(iii) If $f$ is coercive, so is $\tilde{f}$.*

*Proof.* Recall that $\mathcal{F} : \mathbb{R}^n \to \mathbb{S}^m$ by our standing notation. Now let $\mathcal{A} : \mathbb{R}^n \to \mathbb{S}^p$ be *any* affine matrix function such that $\lambda_1 \circ \mathcal{A}$ *is* coercive, and consider the following augmented functions

$$(3.10) \qquad \tilde{\mathcal{F}}_\nu(x) = \begin{pmatrix} \mathcal{A}(x) - \nu I_p & 0 \\ 0 & \mathcal{F}(x) \end{pmatrix} \in \mathbb{S}^{m+p}$$

with $\nu \in \mathbb{N}$ a parameter to be chosen. Define open sets $\Omega_\nu := \{x \in \mathbb{R}^n : \lambda_1(\mathcal{A}(x)) - \nu < \lambda_1(\mathcal{F}(x))\}$, then the entire sequence satisfies $\cup_{\nu=1}^\infty \Omega_\nu = \mathbb{R}^n$. For $\nu$ large enough, let's say for $\nu \geq \nu_0$, the ball $\|x\| \leq R$ is contained in $\Omega_\nu$. By construction

$$\lambda_1\left(\tilde{\mathcal{F}}_\nu(x)\right) = \lambda_1(\mathcal{F}(x)) \text{ for every } x \in \Omega_\nu,$$

16

hence every $f_\nu = \lambda_1 \circ \tilde{\mathcal{F}}_\nu$ with $\nu \geq \nu_0$ is now running as a candidate to becoming the function $\tilde{f}$ we are looking for.

Indeed, consider the linearization of $f_\nu = \lambda_1 \circ \tilde{\mathcal{F}}_\nu$ about some fixed $x$. This is

$$e \mapsto \lambda_1 \left( \tilde{\mathcal{F}}_\nu(x) + \tilde{\mathcal{F}}'_\nu(x) e \right) = \max\{\lambda_1 \left(\mathcal{F}(x) + \mathcal{F}'(x) e\right) ; \lambda_1 \left(\mathcal{A}(x + e)\right) - \nu\},$$

because $\mathcal{A}$ is its own linearization. As $\mathcal{A}$ is coercive, so is $e \mapsto \mathcal{A}(x + e) - \nu I_p$, hence the linearization of each $\tilde{\mathcal{F}}_\nu$ is now coercive, i.e., $f_\nu$ is linearly coercive.

Let $\nu \geq \max\{\nu_0, \epsilon_0\}$, then for every $0 \leq \epsilon \leq \epsilon_0$ and $\|x\| \leq R$ we have

$$\lambda_1 \left(\mathcal{A}(x)\right) - 2\nu \leq \lambda_1 \left(\mathcal{A}(x)\right) - \nu - \epsilon < \lambda_1 \left(\mathcal{F}(x)\right) - \epsilon < \lambda_{r(\epsilon)} \left(\mathcal{F}(x)\right)$$

by the definition of $r(\epsilon)$ and since $x \in \Omega_\nu$. Therefore the first $r(\epsilon)$ eigenvalues of $\tilde{\mathcal{F}}_{2\nu}(x)$ and of $\mathcal{F}(x)$ are the same. The corresponding $r(\epsilon) \times (m + p)$ matrix $\tilde{Q}_\epsilon$ is simply $\tilde{Q}_\epsilon = [0_{p \times r(\epsilon)}; Q_\epsilon^\top]^\top$.

The conclusion is that for $0 \leq \epsilon \leq \epsilon_0$, the $\epsilon$-enlarged subdifferential of $f_{2\nu} = \lambda_1 \circ \tilde{\mathcal{F}}_{2\nu}$ at $x \in \Omega_\nu$ is the same as that of $f = \lambda_1 \circ \mathcal{F}$. In consequence, also steepest $\epsilon$-enlarged descent directions at $x \in \Omega_\nu$ are the same for both representations of $f$. $\square$

This is in contrast with the $\eta$-subdifferentials $\partial_\eta f$ and $\partial_\eta f_{2\nu}$, which are *global* concepts in the sense that they change even when the function is modified far away from the current position $x$. The conclusion is that changing $\partial_\eta f(x)$ into $\partial (f_{2\nu})_\eta (x)$ without affecting $\delta_\epsilon f(x)$ allows to bound the abscissae $t_\eta$. The situations is explained by the following

EXAMPLE. Consider $\mathcal{B}(x) = x^2 \in \mathbb{S}^1$, then $f(x) = \lambda_1 \left(\mathcal{B}(x)\right) = x^2$ is coercive, but the linearizations $e \to \lambda_1 \left(\mathcal{B}(x) + \mathcal{B}'(x)e\right) = x^2 + 2ex$ are not. Now consider the following construction. Let $\tilde{\mathcal{B}}(x) = \text{diag}(-1 + x, -1 - x, x^2) \in \mathbb{S}^3$, then $\lambda_1 \left(\tilde{\mathcal{B}}(x)\right) = x^2$, so $f$ is represented as $\lambda_1 \circ \tilde{\mathcal{B}}$, now on $\mathbb{S}^3$. In the terminology above we have chosen $\mathcal{A}(x) = \text{diag}(x, -x)$ and $\nu = 1$ with $\Omega_\nu = \mathbb{R}$. The linearization of $\tilde{\mathcal{B}}$ about a point $x$ is now $\tilde{\mathcal{B}}(x) + \tilde{\mathcal{B}}'(x)e = \text{diag}(-1 + x + e, -1 - x - e, x^2 + 2ex)$, so $\lambda_1 \left(\tilde{\mathcal{B}}(x) + \tilde{\mathcal{B}}'(x)e\right) = \max\{-1 + |x + e|, x^2 + 2xe\} \to \infty$ as $|e| \to \infty$. In other words, the representation $f = \lambda_1 \circ \tilde{\mathcal{B}}$ is now linearly coercive. $\square$

The first-order bundle algorithm for (1.3) with general $\mathcal{C}^2$ operators is now obtained as follows. We follow the same steps, but use a linearly coercive representation $\tilde{f}$ of $f$ on the compact set $\{x : f(x) \leq f(x_0)\}$ in step 2 and during the search for $t_\eta$ in step 6. The rest of the procedure remains unchanged, and the convergence or finite termination properties are the same.

**3.9. Comments and extensions.** ¿From a practical point of view it may be attractive to modify the bundle method by allowing larger sets of subgradients $\delta_{\epsilon_k} f(x_k) \subset \mathcal{G}_k \subset \partial_{\epsilon_k} f(x_k)$, as proposed in [26]. For instance, some of the elements $g_{k-j} = \mathcal{F}'(x_{k-j})^* G_{k-j}$ from previous steps may be recycled at the step $x_k$. As opposed to the convex case [26], there are two ways how this could be arranged. We could either keep the old $g_{k-j}$, or we could keep only the old $G_{k-j}$ and create new $g_{k-j}^\sharp = \mathcal{F}'(x_k)^* G_{k-j}$ at the actual point $x_k$. We would then accept as trial subgradients $g$ any convex combination $g = \sum_j \alpha_j g_{k-j} + \alpha g'$, $\alpha_j, \alpha \geq 0$, $\sum_j \alpha_j + \alpha = 1$, with $g' \in \delta_{\epsilon_k} f(x_k)$, such that $g \in \partial_{\epsilon_k} f(x_k)$. At any rate, the sets $\mathcal{G}_k$ so obtained are finite extensions of $\delta_{\epsilon_k} f(x_k)$. Notice that our convergence theory covers this case as well, even though the estimates based on $\delta_{\epsilon_k} f(x_k)$ get the more conservative, the larger the gap between $\delta_{\epsilon_k} f(x_k)$ and $\mathcal{G}_k$. Moreover, in each case we have to specify in which way the minimum norm element analogous to (2.7) is computed.

A question of practical importance is whether we can expect a decrease of the order $\mathcal{O}(\eta)$ as in the convex case (see [58]), or whether we will frequently have to be content with the order $\mathcal{O}(\eta^2)$. The following result gives some indication. Concerning terminology, recall that if the maximum eigenvalue $\lambda_1(X)$ has multiplicity $r$, we call $\mathrm{sep}(X) = \lambda_r(X) - \lambda_{r+1}(X) > 0$ the separation of $X$.

LEMMA 3.18. *Let $R > 0$ be fixed. There exist constants $K > 0$ and $\alpha > 0$ such that for every $x$ with $\|x\| \leq R$, every $d$ with $\|d\| = 1$ and every $\eta > 0$ with $f'_\eta(x;d) = \lambda'_1(X + t_\eta D; D)$, the following expansion is valid*

$$\frac{\eta}{t_\eta^2} = \frac{1}{2}\lambda''_1(X;D) + \kappa\, t_\eta$$

*for some $|\kappa| \leq K$ and all $0 < t_\eta \leq \alpha\, \mathrm{sep}(X)$.*

*Proof.* According to Torki [68, Thm. 1.5], the maximum eigenvalue function is twice directionally differentiable and admits an expansion

$$\lambda_1(X + tD) = \lambda_1(X) + t\lambda'_1(X;D) + \frac{t^2}{2}\lambda''_1(X;D) + \mathcal{O}(t^3).$$

Inspecting [66, Thm. 2.8] on which the result is built shows that the $\mathcal{O}(t^3)$ term may more accurately be written as $\mathcal{O}(t^3) = \kappa_1 t^3$, where $|\kappa_1| \leq K_1$ for a constant $K_1$ which is uniform for a bounded set of $X$ and $D$, and for $0 \leq t \leq \alpha\, \mathrm{sep}(X)$, where $\alpha$ is independent of $t$.

A similar directional expansion holds at the second order level. We have

$$\lambda'_1(X + tD; D) - \lambda'_1(X;D) = t\, \lambda''_1(X;D) + \kappa_2 t^2$$

where $|\kappa_2| \leq K_2$ with $K_2$ uniform on a bounded set of $X$ and $D$, but for $t$ only in $0 \leq t \leq \alpha\, \mathrm{sep}(X)$.

Substituting these two estimates with $t = t_\eta$ into (3.5) gives the relationship

$$\eta = \frac{1}{2}t_\eta^2 \lambda''_1(X;D) + \kappa t_\eta^3$$

for $|\kappa| \leq K$ and some constant $K$ which is the same for a bounded set of $X$ and $D$, and for all $t_\eta \leq \alpha\, \mathrm{sep}(X)$. $\square$

Let us substitute this estimate into (3.9) and check whether $t_\eta < \vartheta$, i.e., whether a decrease of the order $\mathcal{O}(\eta)$ may be expected. Taking squares this is equivalent to

$$t_\eta^2 \leq \frac{|f'_\eta(x;d)|\left(|f'_\eta(x;d)| + \sqrt{f'_\eta(x;d)^2 + 4L\rho_0\eta}\right)}{2L^2} + \frac{\rho_0\lambda''_1(X;D)}{2L}t_\eta^2 + Kt_\eta^3$$

where the first term on the right hand side is positive, and the third term is negligeable. Concerning the second term, Torki shows,

$$\lambda''_1(X;D) = \lambda_1\left(2U^\top Q^\top D\left(\lambda_1(X)I_r - X\right)^\dagger DQU\right),$$

where $Q$ is a $r \times m$ matrix whose columns span the eigenspace of $\lambda_1(X)$, and $U$ is a similar matrix for the eigenspace of $\lambda_1(Q^\top DQ)$. This means that $\lambda''_1(X;D)$ behaves roughly like $\mathrm{sep}(X)^{-1} = (\lambda_r(X) - \lambda_{r+1}(X))^{-1}$, the order of magnitude of the pseudo-inverse. This term is expected to explode as the iterates $X_k$ approach a limit $\bar{X}$ with $\lambda_1(\bar{X})$ of multiplicity greater than 1. At least this will happen when the

18

$\lambda_1(X_k)$ have smaller multiplicity than $\lambda_1(\bar{X})$, as is usually the case. In other terms, $\rho_0 \lambda_1''(X; D)/2L$ is expected to be (way) larger than 1, so $t_\eta \leq \vartheta$ is expected to be satisfied asymptotically.

REMARK. Unfortunately this argument remains heuristic, since Torki's expansion at each step $k$ is only valid on a neighborhood $t \leq \alpha \operatorname{sep}(X_k)$ (for the same fixed $\alpha > 0$ which depends neither on $k$, nor on $t$). But it is to be expected that $\operatorname{sep}(X_k) \to 0$, so we cannot a priori render these estimates uniform over $k$. In fact, we could do so provided $\operatorname{sep}(X_k) \to 0$ converged slower than $\|X_k - \bar{X}\| \to 0$. This meant that the $X_k$ approached $\bar{X}$ transversally to the smooth manifold $\mathcal{M}_r = \{X \in \mathbb{S}^m : \lambda_1(X) = \cdots = \lambda_r(X) > \lambda_{r+1}(X)\}$. Only iterates $X_k$ approaching $\mathcal{M}_r$ tangentially cause problems. Ironically, the situation is also o.k. for iterates which are *exactly on* the manifold $\mathcal{M}_r$. Then $\operatorname{sep}(X_k)$ stays away from 0 and Torki's local estimates again hold uniformly over $k$. This strange behavior highlights the non-smooth character of the maximum eigenvalue function. In the past, similar phenomena have motivated approaches, where in order to avoid the tangential zone, iterates have been *forced* to lie on the manifold $\mathcal{M}_r$. We will re-examine this idea in part 2 [49] of this paper.

**3.10. Constrained program.** Let us now address the constrained eigenvalue program (1.4). A first idea, often used in non-smooth optimization, is exact penalization. While this has been reported to induce irregular numerical behavior of bundle methods due to inconveniently large penalty constants (cf. [36]), we believe that the situation is less dramatic for (1.4) with its single scalar constraint. We have the following

PROPOSITION 3.19. *Let $\bar{x}$ be a local minimum of (1.4) such that $\bar{x}$ is not a critical point of $f = \lambda_1 \circ \mathcal{F}$ alone. Then $\bar{x}$ is a KKT-point of (1.4). If at least one of the associated Lagrange multipliers $\bar{\rho} \geq 0$ is known to satisfy $\bar{\rho} \leq \beta$, then $\bar{x}$ is a critical point of the following unconstrained program of the form (1.3):*

$$\min\{c^\top x + \beta \, \lambda_1 \left(\operatorname{diag}\left(0_{1 \times 1}, \mathcal{F}(x)\right)\right) : x \in \mathbb{R}^n\}.$$

*Proof.* Since $c^\top x$ and $f = \lambda_1 \circ \mathcal{F}$ are locally Lipschitz functions, the F. John necessary optimality conditions are satisfied at $\bar{x}$ (see [12, Thm. 6.1.1]). That is, there exist $\bar{\sigma} \geq 0$, $\bar{\rho} \geq 0$, not both zero, such that

$$\bar{\sigma} \, c + \bar{\rho} \, \mathcal{F}'(\bar{x})^* \bar{G} = 0, \ \bar{G} \in \partial\lambda_1\left(\mathcal{F}(\bar{x})\right), \ \bar{\rho} \, \lambda_1\left(\mathcal{F}(\bar{x})\right) = 0, \ \lambda_1\left(\mathcal{F}(\bar{x})\right) \leq 0.$$

Clearly $\bar{\rho} = 0$ is impossible, while $\bar{\sigma} = 0$ would imply that $\bar{x}$ was a KKT point for $\lambda_1 \circ \mathcal{F}$ alone, which was excluded by hypothesis. Therefore $\bar{x}$ is a KKT-point for (1.4). In other terms, we may assume $\bar{\sigma} = 1$, $\bar{\rho} > 0$ above.

Now we show that the set of $\bar{\rho}$ above is bounded. Indeed, suppose on the contrary that we have KKT-conditions with $\bar{\sigma} = 1$, $\rho_n \to \infty$ and $G_n \in \partial\lambda_1\left(\mathcal{F}(\bar{x})\right)$. By compactness of the subdifferential, we may assume $G_n \to G_\infty \in \partial\lambda_1\left(\mathcal{F}(\bar{x})\right)$ for a subsequence. Dividing by $\rho_n$ and passing to the limit then implies $\mathcal{F}'(\bar{x})^* G_\infty = 0$, which means that $\bar{x}$ is a KKT-point for $\lambda_1 \circ \mathcal{F}$ alone, contradicting our hypothesis. Hence the set of $\bar{\rho}$ is bounded.

What we have shown is that program (1.4) is calm at $\bar{x}$ in the sense of Clarke [12, 6.4]. Hence it may be solved by exact penalization. That is, we find $\beta > 0$ such that

$$\min c^\top x + \beta \, \max\{0, \lambda_1\left(\mathcal{F}(\bar{x})\right)\}$$

19

is equivalent to (1.4). But observe that $\max\{0, \lambda_1(\mathcal{F}(\bar{x}))\} = \lambda_1(\text{diag}(0_{1\times 1}, \mathcal{F}(x)))$, so the exact penalty program is of the form (1.3). The fact that every $\beta \geq \bar{\rho}$ will do is standard. $\square$

Let us now look at a second way to address (1.4), which builds on Kiwiel's improvement function [34]. In the convex case $f = \lambda_1 \circ \mathcal{A}$, this has more recently been used by Miller et al. [46, 47], where ideas from [27] have been amalgamated with those of [34]. The emerging numerical method is reported to perform nicely.

Given the current iterate $x_k$ in a minimization algorithm for (1.4), consider the improvement function

$$(3.11)\quad \phi(x, x_k) = \max\left\{c^\top(x - x_k), \lambda_1(\mathcal{F}(x))\right\} = \lambda_1\left(\begin{bmatrix} c^\top(x - x_k) & 0 \\ 0 & \mathcal{F}(x) \end{bmatrix}\right).$$

The following algorithm essentially follows the line of the unconstrained case, where at each instance $k$ the new search direction is computed with respect to the improvement function $\phi(x_k; \cdot)$ instead of $f$.

<div align="center">Spectral bundle algorithm for program (1.4)</div>

| | |
|---|---|
| 1. | Let $\omega, \rho_0, \theta_0, \epsilon_0^\sharp, x_0, \gamma_k, S, F \subset \mathbb{N}$ and `slope` $\in \{\texttt{steep}, \texttt{flat}\}$ be as in the unconstrained algorithm. |
| 2. | Given the current $x_k$, stop if $0 \in \partial\phi(x_k; x_k)$. Otherwise let $\bar{\epsilon}_k$ such that $0 \in \delta_{\bar{\epsilon}_k}\phi(x_k; x_k)$, but $0 \notin \delta_\epsilon\phi(x_k; x_k)$ for $\epsilon < \bar{\epsilon}_k$. Choose $\epsilon_k \leq \epsilon_k^\sharp$ such that $\Delta_{\epsilon_k}(X_k) \geq \min\{\bar{\epsilon}_k/m, \epsilon_k^\sharp/m\}$. |
| 3. | compute a direction $d_k$ of approximate steepest $\epsilon_k$-enlarged descent at $x_k$ with respect to the function $\phi(x_k; \cdot)$. |
| 4. | If $|\tilde{\phi}'(x_k; \cdot)_{\epsilon_k}(x_k; d_k)| \leq \epsilon_k^\sharp$ put `slope` $=$`flat`, otherwise put `slope` $=$ `steep`. |
| 5. | compute $\eta_k$ as in (3.2) using $\phi(x_k; \cdot)$. |
| 6. | Search for $t_{\eta_k}$ as in the unconstrained algorithm using $\phi(x_k; \cdot)$ and the corresponding $\mathcal{F}_k = \text{diag}\{c^\top(x - x_k), \mathcal{F}\}$ in lieu of $f$ and $\mathcal{F}$. |
| 7.-9. | In analogy with the unconstrained case. |

Then we have the following

THEOREM 3.20. *Consider program (1.4). Let $x_0$ be fixed and suppose $\mathcal{F}$ is of class $\mathcal{C}_b^2$. Suppose $c^\top x$ is bounded below on the feasible set $\{x \in \mathbb{R}^n : f(x) \leq 0\}$. Let the sequence $x_k$ starting with $x_0$ be generated by the spectral bundle algorithm for (1.4). Then the following cases may occur:*
  1. *All iterates are infeasible, i.e., $f(x_k) > 0$ for all $k$. If $\bar{x}$ is an accumulation point of the entire sequence $x_k$ (when $F$ is finite) and of the flat subsequence $x_k$, $k \in F$ (when $F$ is infinite), then $\bar{x}$ is a critical point of $f$.*
  2. *The iterates $x_k$ are strictly feasible, i.e., $f(x_k) < 0$, for some $k_0$ and all $k \geq k_0$. If $\bar{x}$ is an accumulation point of the entire sequence $x_k$ (if $F$ is finite) and of the flat subsequence $x_k$, $k \in F$ otherwise, then $\bar{x}$ satisfies the F. John necessary optimality conditions for program (1.4).*

*Proof.* 1) Let us first examine the situation where the initial point $x_0$ is infeasible, $f(x_0) > 0$. Then there are two possibilities. Either feasibility is reached in finite time, i.e., $f(x_k) \leq 0$ at some stage $k$. Or $f(x_k) > 0$ for all $k$, so that feasibility is never reached. In the second case $\phi(x_k, x) = f(x)$ around $x_k$, and the algorithm essentially

<div align="center">20</div>

behaves like the unconstrained bundle algorithm, that is, it reduces $f$. Since the $f(x_k)$ are bounded below by 0, the same conclusions are obtained. More precisely, every accumulation point $\bar{x}$ of $x_k$ is a critical point of $f$ if the set $F$ is finite, and so is every accumulation point of the flat subsequence if $F$ is infinite. Such an accumulation point may or may not be feasible.

2) Let us now suppose that some iterate $k$ is feasible, $f(x_k) \leq 0$. Then $\phi(x_k; x_k) = 0$. Unless the algorithm halts with $0 \in \partial \phi(x_k; x_k)$, the new $d_k$ is a direction of descent of $\phi(x_k; \cdot)$ at $x_k$, so the line search will give a descent step with $\phi(x_k; x_{k+1}) < \phi(x_k; x_k) = 0$. Then $f(x_{k+1}) \leq \phi(x_k; x_{k+1}) < 0$. Therefore iterate $x_{k+1}$ is even strictly feasible.

3) Suppose next that $f(x_k) < 0$ for some $k$. Then by repeating the argument in 2), all iterates $x_j$, $j \geq k$ stay strictly feasible. What is more, $c^\top (x_{k+1} - x_k)) \leq \phi(x_k; x_{k+1}) < \phi(x_k; x_k) = 0$, so the algorithm now reduces the objective function at each step and continues to do so during the following steps.

Altogether iterates now decrease in value and remain strictly feasible. Since by hypothesis the problem is bounded below, the series $\sum_k \alpha_{t_k}$ converges. By construction, this yields then the same cases as discussed in the unconstrained algorithm.

4) Suppose for instance that some subsequence $k \in \mathcal{N}$ has $\bar{\epsilon}_k \to 0$, where $0 \in \delta_{\bar{\epsilon}_k} \phi(x_k; x_k)$. Let $\bar{x}$ be one of its accumulation points, then the argument of Lemma 3.2 shows $0 \in \partial \phi(\bar{x}; \bar{x})$. In that case $\bar{x}$ satisfies the F. John necessary optimality conditions, so it must be either a KKT-point, or a critical point of $f$ alone.

Feasibility $f(\bar{x}) \leq 0$ is clear. But $f(\bar{x}) < 0$ is not possible, for in that case we would have $\partial \phi(\bar{x}; \bar{x}) = \{c\}$, a contradiction. So $f(\bar{x}) = 0$. Then both $c^\top (x - x_k)$ and $\mathcal{F}(x)$ are active at $x_k$, hence $0 = \alpha c + (1 - \alpha) g$ for some $0 \leq \alpha \leq 1$ and $g \in \partial f(\bar{x})$. If $\alpha > 0$ we have a KKT point. If $\alpha = 0$, then $\bar{x}$ is a critical point of $f$ alone, whose value is 0.

5) Suppose next that $\tilde{\phi}'_{\epsilon_k}(x_k; \cdot)(x_k; d_k) \to 0$ in tandem with $\epsilon_k \to 0$ for a subsequence $k \in \mathcal{N}$. Here the argument of Lemma 3.2 shows that every accumulation point $\bar{x}$ of $x_k$, $k \in \mathcal{N}$, has $0 \in \partial \phi(\bar{x}; \bar{x})$, so the conclusion is the same.

6) Finally suppose $0 \in \partial \phi(x_k; x_k)$ at some $k$, in which case the algorithm halts. Here if $f(x_k) > 0$, we must have a critical point of $f$ alone. If $f(x_k) \leq 0$, the discussion is the same as in 4) above. This completes the proof. $\square$

REMARK. Clearly when $0 \in \partial f(\bar{x})$ with $f(\bar{x}) > 0$, the algorithm fails. While this always happens when the problem is infeasible, it is clear that even in the feasible case we may create situations, where a first order method like the proposed one must fail. For instance, we could arrange that feasible points can only be reached from $x_0$ by *increasing* the improvement function, which the method never does. Nonetheless, Theorem 3.20 seems practically useful, as local minima or critical points of $f$ alone are not expected to occur frequently in practice.

Notice that for a convex constraint, $f(x) = \lambda_1 (\mathcal{A}(x)) \leq 0$, the algorithm never fails. We have the following

COROLLARY 3.21. *Suppose $f = \lambda_1 \circ \mathcal{A}$ is convex and program (1.4) is bounded below and has a strictly feasible point. Then every accumulation point of the sequence $x_k$ generated by the constraint bundle algorithm is a minimum of (1.4).*

*Proof.* Suppose we had $f(x_k) > 0$ for all $k$. Then some accumulation point $\bar{x}$ of the $x_k$ is critical for $f$ alone, which by convexity means $\bar{x}$ is a minimum of $f$. Suppose $f(\bar{x}) > 0$, then the program has no feasible points, a contradiction. But $f(\bar{x}) = 0$ is also impossible, because no point is strictly feasible. This means $f(x_k)$ become strictly feasible after a finite number of iterations. Now Theorem 3.20 shows that

every accumulation point $\bar{x}$ of the $x_k$ satisfies the F. John optimality condition. If $\bar{x}$ is a KKT-point, then it is a minimum by convexity, so we are done. The other possibility is that $\bar{x}$ is a minimum of $f$. But the proof of Theorem 3.20 shows that $f(\bar{x}) = 0$, while a minimum of $f$ (if any) must have strictly negative value. So this is impossible. This completes the proof. □

**4. Conclusion.** We have proved a global convergence result for constrained and unconstrained optimization problem using non-convex maximum eigenvalue functions, which assures subsequence convergence of the sequence of iterates towards critical points under mild assumptions. The proposed method computes qualified descent steps using an inner approximation of $\epsilon$-subdifferentials proposed in [13] and [58]. Recent numerical tests, to be published in [3, 4], seem to indicate that the method performs fairly well in practice. An extension to second order methods will be presented in part 2 [49] of this work.

REFERENCES

[1] P. APKARIAN, D. NOLL, H.D. TUAN, *A prototype primal-dual LMI interior point algorithm for non-convex robust control problems*, Rapport interne 01-08 (2002), UPS-MIP, CNRS UMR 5640. `http://mip.ups-tlse.fr/publi/publi.html`

[2] P. APKARIAN, D. NOLL, H.D. TUAN, *Fixed-order $\mathcal{H}_\infty$ control design via an augmented Lagrangian method*, Int. J. Robust and Nonlin. Control, vol. 13, no. 12, 2003, pp. 1137 – 1148.

[3] P. APKARIAN, D. NOLL, *Controller design via nonsmooth multi-directional search*. Submitted.

[4] P. APKARIAN, D. NOLL, *Nonsmooth $H_\infty$-synthesis*. Submitted.

[5] P. APKARIAN, H.D. TUAN, *Concave programming in control theory*, J. Global Optim., 15 (1999), pp. 343 – 370.

[6] P. APKARIAN, H.D. TUAN, *Robust control via concave minimization - local and global algorithms*, IEEE Trans. on Autom. Control, 45 (2000), pp. 299 – 305.

[7] V. I. ARNOLD, *On matrices depending on parameters*, Russ. Math. Surveys, 26 (1971), pp. 29 – 43.

[8] S. BOYD, L. ELGHAOUI, E. FERON, V. BALAKRISHNAN, *Linear Matrix Inequalities in System and Control Theory*, vol. 15 of *Studies in Appl. Math.* SIAM, Philadelphia, 1994.

[9] E.B. BERAN, L. VANDENBERGHE, S. BOYD, *A global BMI algorithm based on the generalized Benders decomposition*, Proc. European Control Conf., Belgium, 1997.

[10] J. F. BONNANS, *Local analysis of Newton-type methods for variational inequalities and nonlinear programming*, Appl. Math. Opt., 29 (1994), pp. 161 – 186.

[11] S. BOYD, L. VANDENBERGHE, *Semidefinite programming relaxations of non-convex problems in control and combinatorial optimization*.

[12] F.H. CLARKE, *Optimization and Nonsmooth Analysis*, John Wiley, New York, 1983.

[13] J. CULLUM, W.E. DONATH, P. WOLFE, *The minimization of certain nondifferential sums of eigenvalues of symmetric matrices*, Math. Progr. Stud., vol. 3 (1975), pp. 35 – 55.

[14] J.E. DENNIS JUN., R.B. SCHNABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice Hall Series in Computational Mathematics, 1983.

[15] L. ELGHAOUI, V. BALAKRISHNAN, *Synthesis of fixed-structure controllers via numerical optimization*, Proc. IEEE Conf. Dec. Contr., December 1994.

[16] B. FARES, D. NOLL, P. APKARIAN, *An augmented Lagrangian method for a class of LMI-constrained problems in robust control theory*, Intern. J. Control, 74 (2001), pp. 348 – 360.

[17] B. FARES, D. NOLL, P. APKARIAN, *Robust control via successive semidefinite programming*, SIAM J. on Control Optim., 40 (2002), pp. 1791 – 1820.

[18] R. FLETCHER, *Semidefinite constraints in optimization*, SIAM J. Control Optim., 23 (1985), pp. 493 – 513.

[19] R. FLETCHER, S. LEYFFER, *A bundle filter method for nonsmooth nonlinear optimization*. University of Dundee, Report NA/195, 1999.

[20] A. FUDULI, M. GAUDIOSO, G. GIALLOMBARDO, *Minimizing nonconvex nonsmooth functions via cutting planes and proximity control*, SIAM J. on Optim., to appear.

[21] M. FUKUDA, M. KOJIMA, *Branch-and-cut algorithms for the bilinear matrix inequality eigenvalue problem*, Research Report on Mathematical and Computational Sciences. Series B:

Operations Research, 1999.

[22] K.M. GRIGORIADIS, R.E. SKELTON, *Low-order control design for LMI problems using alternating projection methods*, Automatica, 32 (1996), pp. 1117 – 1125.

[23] K.-C. GOH, M.G. SAFONOV, G.P. PAPAVASSILOPOULOS, *Global optimization for the biaffine matrix inequality problem*, J. Global Optim., 7 (1995), pp. 365 – 380.

[24] C. HELMBERG, *A cutting plane algorithm for large scale semidefinite relaxations.* ZIB-Report ZR-01-26. To appear in Padberg Festschrift *The sharpest cut*, SIAM, 2001.

[25] C. HELMBERG, K.C. KIWIEL, *A spectral bundle method with bounds*, Math. Programming, 93, no. 2, 2002, pp. 173 – 194.

[26] C. HELMBERG, F. OUSTRY,*Bundle methods to minimize the maximum eigenvalue function,* In: L. VANDENBERGHE, R. SAIGAL, H. WOLKOWITZ, (eds.), Handbook of Semidefinite Programming. Theory, Algorithms and Applications, vol. 27, Int. Series in Oper. Res. and Management Sci. Kluwer Acad. Publ., 2000 .

[27] C. HELMBERG, F. RENDL,*A spectral bundle method for semidefinite programming,* SIAM J. Optim., 10, no. 3 (2000), pp. 673 – 696.

[28] C. HELMBERG, F. RENDL,*Solving quadratic $(0,1)$-problems by semidefinite programs and cutting planes*, Math. Programming, 82 (1998), pp. 291 – 315.

[29] J.W. HELTON, O. MERINO, *Coordinate optimization for bi-convex matrix inequalities*, Proc. IEEE Conf. on Decision and Control, San Diego, 1997, pp. 3609 – 3613.

[30] J.-B. HIRIART-URRUTY, C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms, part II,* Springer Verlag, Berlin, 1993.

[31] C.W.J. HOL, C.W. SCHERER, E.G. VAN DER MECHÉ, O.H. BOSGRA, *A nonlinear SDP approach to fixed-order controller synthesis and comparison with two other methods applied to an active suspension system*, submitted, 2002.

[32] T. IWASAKI, *The dual iteration for fixed order control*, Proc. Amer. Control Conf., 1997, pp. 62 – 66.

[33] T. KATO, *Perturbation Theory for Linear Operators.* Springer Verlag, 1984.

[34] K.C. KIWIEL, *Methods of descent for nondifferentiable optimization*, Lect. Notes in Math. vol. 1133, Springer Verlag, Berlin, 1985.

[35] K.C. KIWIEL, *A linearization algorithm for optimizing control systems subject to singular value inequalities*, IEEE Trans. Autom. Control AC-31, 1986, pp. 595 – 602.

[36] K.C. KIWIEL, *A constraint linearization method for nondifferentiable convex minimization*, Numerische Mathematik 51, 1987, pp. 395 – 414.

[37] K.C. KIWIEL, *Proximity control in bundle methods for convex nondifferentiable optimization*, Math. Programming 46, 1990, pp. 105 – 122. Math. Programming, 46 (1990), pp. 105 – 122.

[38] K.C. KIWIEL, *Restricted-step and Levenberg-Marquardt techniques in proximal bundle methods for nonconvex nondifferentiable optimization*, SIAM J. on Optim. 6 (1996), pp. 227 – 249.

[39] K. KRISHNAN, J.E. MITCHELL, *Cutting plane methods for semidefinite programming*, submitted, 2002.

[40] C. LEMARÉCHAL, *Extensions diverses des méthodes de gradient et applications*, Thèse d'Etat, Paris, 1980.

[41] C. LEMARÉCHAL, *Bundle methods in nonsmooth optimization.* In: Nonsmooth Optimization, Proc. IIASA Workshop 1977, C. LEMARÉCHAL, R. MIFFLIN (eds.), Pergamon Press, 1978.

[42] C. LEMARÉCHAL, *Nondifferentiable Optimization*, chapter VII in: *Handbooks in Operations Research and Management Science,* vol. 1, Optimization, G.L. Nemhauser, A.H.G. Rinnooy Kan, M.J. Todd (eds.), North Holland, 1989.

[43] C. LEMARÉCHAL, A. NEMIROVSKII, Y. NESTEROV, *New variants of bundle methods,* Math. Programming, 69 (1995), pp. 111 – 147.

[44] C. LEMARÉCHAL, F. OUSTRY,*Nonsmooth algorithms to solve semidefinite programs.* In: L. EL GHAOUI, L. NICULESCU (eds.), Advances in LMI methods in Control, SIAM Advances in design and Control Series, 2000.

[45] C. LEMARÉCHAL, F. OUSTRY, C. SAGASTIZABAL, *The $\mathcal{U}$-Lagrangian of a convex function*, Trans. Amer. Math. Soc. 352, 2000.

[46] S.A. MILLER, R.S. SMITH, *Solving large structured semidefinite programs using an inexact spectral bundle method*, Proc. IEEE Conf. Dec. Control, 2000, pp. 5027 – 5032.

[47] S.A. MILLER, R.S. SMITH, *A bundle method for efficiently solving large structured linear matrix inequalities.* Proc. Amer. Control Conf., Chicago, 2000.

[48] A. NEMIROVSKII, *Several NP-Hard problems arising in robust stability analysis,* Mathematics of Control, Signals, and Systems, vol. 6, no. 1 (1994), pp. 99 – 105.

[49] D. NOLL, P. APKARIAN, *Spectral bundle methods for nonconvex maximum eigenvalue functions. Part 2: second-order methods.* Mathematical Programming, series B, to appear.

[50] D. NOLL, M. TORKI, P. APKARIAN, *Partially augmented Lagrangian method for matrix inequalities.* Submitted.

[51] M.L. OVERTON, *On minimizing the maximum eigenvalue of a symmetric matrix,* SIAM J. Matrix Anal. Appl., vol. 9, no. 2 (1988), pp. 256 − 268.

[52] M.L. OVERTON, *Large-scale optimization of eigenvalues,* SIAM J. Optim., vol. 2, no. 1 (1992), pp. 88 − 120.

[53] M.L. OVERTON, R.S. WOMERSLEY,*Optimality conditions and duality theory for minimizing sums of the largest eigenvalues of symmetric matrices,* Math. Programming, 62 (1993), pp. 321 − 357.

[54] M.L. OVERTON, R.S. WOMERSLEY, *Second derivatives for optimizing eigenvalues of symmetric matrices,* SIAM J. Matrix Anal. Appl., vol. 16, no. 3 (1995), pp. 697 − 718.

[55] J.-P. A. HAEBERLY, M.L. OVERTON,*A hybrid algorithm for optimizing eigenvalues of symmetric definite pencils.* SIAM J. Matrix Anal. Appl., 15 (1994), pp. 1141 − 1156.

[56] F. OUSTRY, *Vertical development of a convex function,* J. Convex Analysis, 5 (1998), pp. 153 − 170.

[57] F. OUSTRY, *The $\mathcal{U}$-Lagrangian of the maximum eigenvalue function,* SIAM J. Optim., vol. 9, no. 2 (1999), pp. 526 − 549.

[58] F. OUSTRY, *A second-order bundle method to minimize the maximum eigenvalue function,* Math. Programming Series A, vol. 89, no. 1 (2000), pp. 1 − 33.

[59] M. ROTUNNO, R.A. DE CALLAFON, *A bundle method for solving the fixed order control problem,* Proc. IEEE Conf. Dec. Control, Las Vegas, 2002, pp. 3156 − 3161.

[60] C. SCHERER, *A full block S-procedure with applications,* Proc. IEEE Conf. on decision and Control, San Diego, USA, 1997, pp. 2602 − 2607.

[61] C. SCHERER, *Robust mixed control and linear parameter-varying control with full block scalings,* SIAM Advances in Linear Matrix Inequality Methods in Control Series, L. El Ghaoui, S.I. Niculescu (eds.), 2000.

[62] H. SCHRAMM, J. ZOWE, *A version of the bundle method for minimizing a nonsmooth function: conceptual idea, convergence analysis, numerical results,* SIAM J. Optim., 2 (1992), pp. 121 − 152.

[63] A. SHAPIRO, *Extremal problems on the set of nonnegative matrices,* Lin. Algebra and its Appl., 67 (1985), pp. 7 − 18.

[64] A. SHAPIRO, *First and second order analysis of nonlinear semidefinite programs,* Math. Programming Ser. B, 77 (1997), pp. 301 − 320.

[65] A. SHAPIRO, M.K.H. FAN,*On eigenvalue optimization,* SIAM J. Optim., vol. 5, no. 3 (1995), pp. 552 − 569.

[66] G.W. STEWARD, JI-GUANG SUN, *Matrix Perturbation Theory,* Computer Science and Scientific Computing, Academic Press, 1990.

[67] M. TORKI,*First- and second-order epi-differentiability in eigenvalue optimization,* J. Math. Anal. Appl., 234 (1999), pp. 391 − 416.

[68] M. TORKI, *Second order directional derivative of all eigenvalues of a symmetric matrix,* Nonlin. Analysis, Theory, Methods and Appl., 46 (2001), pp. 1133 − 1500.

[69] H.D. TUAN, P. APKARIAN, Y. NAKASHIMA, *A new Lagrangian dual global optimization algorithm for solving bilinear matrix inequalities,* Proc. Amer. Control Conf., 1999.

[70] J.F. WHIDBORNE, J. WU, R.H. ISTEPANIAN, *Finite word length stability issues in an $\ell_1$ Framework,* Int. J. Control, 73, no. 2 (2000), pp. 166 − 176.

[71] PH. WOLFE, *A method of conjugate subgradients for minimizing nondifferentiable convex functions,* Math. Programming Studies, 3 (1975), pp. 145 − 173.