

# A simple Newton method for local nonsmooth optimization

A.S. Lewis <sup>\*</sup>      C.J.S. Wylie <sup>†</sup>

July 30, 2019

## Abstract

Superlinear convergence has been an elusive goal for black-box nonsmooth optimization. Even in the convex case, the subgradient method is very slow, and while some cutting plane algorithms, including traditional bundle methods, are popular in practice, local convergence is still sluggish. Faster variants depend either on problem structure or on analyses that elide sequences of “null” steps. Motivated by a semi-structured approach to optimization and the sequential quadratic programming philosophy, we describe a new bundle Newton method that incorporates second-order objective information with the usual linear approximation oracle. One representative problem class consists of maxima of several smooth functions, individually inaccessible to the oracle. Given as additional input just the cardinality of the optimal active set, we prove local quadratic convergence. A simple implementation shows promise on more general functions, both convex and nonconvex, and suggests first-order analogues.

**Key words:** nonsmooth optimization, bundle method, Newton method, sequential quadratic programming, nonconvex

**AMS 2010 Subject Classification:** 90C25, 65K05, 49M15

## 1 Introduction

Accurate deterministic minimization of a smooth objective function  $f: \mathbf{R}^n \rightarrow \mathbf{R}$  constitutes the foundation of the classical optimization literature [28]. Typically we assume a black-box oracle, returning, at any point  $x \in \mathbf{R}^n$ , the objective value

---

<sup>\*</sup>ORIE, Cornell University, Ithaca, NY 14853, USA. [people.orie.cornell.edu/aslewis](http://people.orie.cornell.edu/aslewis)  
Research supported in part by National Science Foundation Grant DMS-1613996.

<sup>†</sup>ORIE, Cornell University, Ithaca, NY 14853, USA.

$f(x)$ , gradient  $\nabla f(x)$ , and perhaps the Hessian  $\nabla^2 f(x)$ . Fast local convergence — superlinear or quadratic — relies ultimately on Newton’s method.

By contrast, nonsmooth optimization, even for convex objectives, is more challenging, at least in the analogous black-box model [23]. Armed with the right tools for the circumstances — some manner of global structural knowledge — reasonably accurate nonsmooth convex minimization may be tractable. The contemporary optimization literature highlights various examples: the classical proximal point method and its various splitting-type extensions, when we can suitably decompose the objective into smooth and prox-friendly ingredients [31, 33]; basic techniques for smooth minimization (along the lines surveyed in [25]) extended to nonsmooth objectives amenable to suitable smoothing techniques [26]; polynomial-time interior-point techniques for semidefinite-representable minimization [27].

Outside such structured convex optimization realms, however, we fall back on black-box techniques, using function values and subgradients but blind to any global structure beyond perhaps convexity [2, 25]. The simplicity and all-purpose appeal of such algorithms is tempered by their local convergence in practice. At large scale, the classical subgradient algorithm is the notoriously slow method of last resort [23]. In the convex case, modern cutting plane methods based on analytic centers [1, 7, 24] or the volumetric barrier [11, 36], can be highly effective at moderate scale, but convergence is linear.

Bundle methods, which originate with [12, 19], have proved quite successful in large-scale practice, are supported by convergence theory [6, 9], and extend to the nonconvex case (although implementation is delicate — see [10]). However, classical bundle methods converge sublinearly, and while more recent black-box versions such as [21] are faster, they too are complex to implement, and the convergence analysis considers only “serious” steps rather than oracle calls. Sagastizábal’s 2018 ICM lecture [34] includes an elegant and comprehensive survey. Impractically large subproblems disadvantage some otherwise promising analogous developments, like level bundle methods [13].

The quest for superlinear acceleration described in [34] has for two decades also driven important “ $\mathcal{VU}$ ” and “partly smooth” theory, pioneered by Lemaréchal, Mifflin, Oustry, and Sagastizábal [14, 15, 20] and paralleled in [16]. However, still missing in black-box nonsmooth optimization, nonconvex or convex, is the fast local convergence of Newton’s method. Our aim here is just such a black-box Newton-type local optimization method, supported by a superlinear convergence theory on a range of interesting nonsmooth functions, and promising in practice.

As a step in that direction, we take an unconventional *semi-structured* approach. As motivation, we first consider how to minimize a common type of structured nonsmooth objective function, specifically a pointwise maximum of finitely-many smooth component functions, but using a black-box oracle for the objective that *cannot access the component functions individually*. The resulting algorithm becomes a model-based black-box “bundle Newton” method (terminology coming from [18])

that converges locally quadratically on nonsmooth functions of this max-type, but that also shows promise in practice on more general nonsmooth functions, both convex and nonconvex.

Like any local method, including Newton’s method itself, globalizing the algorithm is an important challenge, and a topic for future work, as is the first-order analogue we describe briefly at the conclusion of this work. Nonetheless, this Newton-type local approach seems a promising new ingredient for nonsmooth optimization.

## 2 An algorithm for nonsmooth minimization

We seek a local method for finding a minimizer  $\bar{x} \in \mathbf{R}^n$  of a continuous nonsmooth objective function  $f$ . We assume, around every point in some set,  $\mathcal{D} \subset \mathbf{R}^n$ , that the objective  $f$  is twice continuously differentiable. Given any point  $s \in \mathcal{D}$ , an oracle returns the value  $f(s)$ , gradient  $\nabla f(s)$ , and Hessian  $\nabla^2 f(s)$ , allowing us to build the corresponding linear and quadratic approximations,  $l_s$  and  $q_s$ , to  $f$ . Beyond that, we call on no further information about the function  $f$ .

Since the objective  $f$  is nonsmooth, we cannot typically hope to find a point  $s \in \mathcal{D}$  where the gradient  $\nabla f(s)$  is small. Instead we seek a finite set (or *bundle*)  $S \subset \mathcal{D}$  with small diameter

$$\text{diam } S = \max\{|s - s'| : s, s' \in S\}$$

(where  $|\cdot|$  denotes the Euclidean norm), and small *optimality measure*

$$\Theta(S) = \min |\text{conv}(\nabla f(S))|.$$

We call the elements of the bundle *reference points*. In terms of the simplex

$$\Delta_S = \left\{ \lambda \in \mathbf{R}_+^S : \sum_{s \in S} \lambda_s = 1 \right\},$$

the optimality measure is

$$(2.1) \quad \Theta(S) = \min_{\lambda \in \Delta_S} \left| \sum_{s \in S} \lambda_s \nabla f(s) \right|.$$

If the set of gradients  $\nabla f(S)$  is affinely independent, then the minimization problem (2.1) has a unique optimal solution  $\lambda$ . This vector, which plays an important role in the algorithm, we call the *Lagrange multiplier estimate*.

We describe an algorithm that iteratively updates a bundle  $S$  one reference point at a time, stopping if it encounters a point outside the set  $\mathcal{D}$ . Thus  $S$  has a fixed size (cardinality) throughout: a judicious choice of this integer parameter, denoted  $k$ , is a crucial feature of the algorithm that distinguishes it from more standard cutting plane and bundle methods. Too small a value of  $k$  causes the method to fail because the optimality measure  $\Theta(S)$  remains bounded away from zero; too large a value may lead to an ill-defined Lagrange multiplier estimate.

## The role of convexity

The core algorithm below is aimed at strongly convex objectives  $f$ , even though its statement makes sense without convexity and our interest in it is purely local. Indeed, much of the analysis that follows is independent of convexity, and a straightforward modification of the core method results in our culminating nonconvex algorithm.

Nonetheless, simple examples show that the unmodified method can fail for nonconvex objectives. We present this algorithm first because it is easier to motivate, our aim in Section 3. Most importantly, the final step in our key quadratic convergence proof (Theorem 5.11) depends crucially on a convexity argument. Without more ado, here is an informal description of the core algorithm.

### Algorithm 2.2 ( $k$ -bundle Newton method for strongly convex $f$ )

**Require:** initial bundle  $S \subset \mathcal{D}$  of size  $k$ , tolerances  $\bar{\epsilon}, \bar{\delta} \geq 0$ ;

**for** iteration = 1, 2, 3, ... **do**

**for**  $s \in S$  **do**

$$l_s(\cdot) = f(s) + \nabla f(s)^T(\cdot - s);$$

$$q_s(\cdot) = l_s(\cdot) + \frac{1}{2}(\cdot - s)^T \nabla^2 f(s)(\cdot - s);$$

**end for**

$$\delta = \Theta(S);$$

  choose  $\lambda \in \Delta_S$  with  $|\sum_{s \in S} \lambda_s \nabla f(s)| = \delta$ ;

**if**  $\text{diam } S < \bar{\epsilon}$  and  $\delta < \bar{\delta}$  **then**

**return** *Stopped: nearly optimal*;

**else**

    choose  $\hat{x} \in \text{argmin} \{ \sum_{s \in S} \lambda_s q_s(x) : x \in \mathbf{R}^n, l_s(x) \text{ equal for all } s \in S \}$

**end if**

**if**  $\hat{x} \notin \mathcal{D}$  **then**

**return** *Stopped: nonsmooth point.*;

**else**

    choose  $s \in S$  minimizing  $\Theta((S \setminus \{s\}) \cup \{\hat{x}\})$ ;

$$S = (S \setminus \{s\}) \cup \{\hat{x}\};$$

**end if**

**end for**

Notice that the case  $k = 1$  is just the classical Newton method.

## 3 Motivating the method

### 3.1 The optimality measure

Basic aspects of the algorithm coincide with traditional cutting plane methods. The method aims to construct a sequence of bundles  $S$  converging to the minimizer  $\bar{x}$

(and hence in particular with diameter converging to zero) in such a way that the optimality measure  $\delta = \Theta(S)$  also converges to zero: in that case we say that the algorithm *succeeds*. In practice we deduce approximate optimality when both these measures are small.

The algorithm computes an optimal vector  $\lambda$  (the Lagrange multiplier estimate) for the optimality measure expression (2.1). If we define the weighted average of the reference points

$$(3.1) \quad \bar{s} = \sum_{s \in S} \lambda_s s,$$

then  $f(\bar{s})$  is certainly an upper bound on the optimal value  $\min f$ . To obtain an approximate lower bound, we use convexity to note that each linear approximation

$$l_s(x) = f(s) + \nabla f(s)^T(x - s) \quad (x \in \mathbf{R}^n),$$

minorizes  $f$ , and hence so does their weighted average  $\sum_s \lambda_s l_s$ . Denoting by  $L$  a Lipschitz constant for  $f$  on some convex set containing all the reference points, we have  $|\nabla f(s)| \leq L$  for all  $s \in S$ . Then for all points  $x \in \mathbf{R}^n$  we have

$$\begin{aligned} f(\bar{s}) - L \text{diam}(S) &\leq f(\bar{s}) - \sum_{s \in S} \lambda_s L |s - \bar{s}| \leq f(\bar{s}) + \sum_{s \in S} \lambda_s \nabla f(s)^T (\bar{s} - s) \\ &= \sum_{s \in S} \lambda_s l_s(\bar{s}) \leq \sum_{s \in S} \lambda_s l_s(x) + \delta |x - \bar{s}| \leq f(x) + \delta |x - \bar{s}|, \end{aligned}$$

so in conclusion we have

$$(3.2) \quad \min f \leq f(\bar{s}) \leq \min\{f + \delta|\cdot - \bar{s}|\} + L \text{diam}(S).$$

Thus if both the diameter of the bundle  $S$  and the optimality measure  $\delta$  are small, then the objective value at the point  $\bar{s}$  lies in the small interval between the minimum values of the objective function and a slightly perturbed function. In this sense, the current bundle constitutes an approximate certificate of optimality.

As a consequence of this argument, success of the algorithm requires a certain lower bound on the bundle size  $k$ , as we discuss next.

### 3.2 A lower bound on bundle size: Carathéodory number

Given any set  $\Gamma \subset \mathbf{R}^n$  containing zero in its convex hull, define the *Carathéodory number*  $\text{car } \Gamma$ , to be the minimum size of a subset whose convex hull contains zero. By Carathéodory's theorem, we see

$$1 \leq \text{car } \Gamma \leq 1 + \dim(\text{conv } \Gamma).$$

Suppose that the set  $\mathcal{D}$  has full measure. When  $\mathcal{D}$  is simply the set of all points where  $f$  is twice continuously differentiable, this assumption typically holds in practice, and in particular if the objective  $f$  is semi-algebraic [5]. Define the limiting gradient set

$$(3.3) \quad \Gamma = \left\{ \lim_r \nabla f(x_r) : \lim_r x_r = \bar{x}, x_r \in \mathcal{D} \text{ for } r = 1, 2, \dots \right\}.$$

Since  $f$  is locally Lipschitz, being continuous and convex, we have (see [4])

$$(3.4) \quad 0 \in \partial f(\bar{x}) = \text{conv } \Gamma.$$

For the algorithm to succeed, the optimality measure  $\Theta(S)$  must converge to zero for some sequence of bundles  $S$  converging to the minimizer  $\bar{x}$ , so zero is a convex combination of  $k$  elements of the set  $\Gamma$ , and hence the lower bound

$$k \geq \text{car } \Gamma$$

must hold.

### 3.3 The active subspace

Having computed the optimality measure  $\delta = \Theta(S)$ , the method next seeks a new reference point. If the current bundle  $S$  is close to the minimizer  $\bar{x}$ , then the cutting plane model  $\tilde{f}: \mathbf{R}^n \rightarrow \mathbf{R}$  defined by

$$\tilde{f}(x) = \max_{s \in S} l_s(x),$$

minorizes the objective  $f$ , and approximates it around  $\bar{x}$ . Furthermore, at every point  $x$  on the *active subspace* where all the linear approximations are equal,

$$M = \{x \in \mathbf{R}^n : l_s(x) \text{ all equal for all } s \in S\},$$

the cutting plane model has subdifferential

$$\partial \tilde{f}(x) = \text{conv}(\nabla f(S))$$

and hence nonsmooth slope (the fastest rate of decrease) equal to  $\delta$ . Thus when the optimality measure is small, the cutting plane model is approximately minimized throughout the active subset, so there is where the method seeks a new reference point.

### 3.4 An upper bound on bundle size: affine independence

For algorithmic stability, the minimizer  $\bar{x}$  should be close to the active subspace  $M$ . Since the bundle  $S$  is close to  $\bar{x}$ , the values of the linear approximations  $l_s(\bar{x})$  are all close to  $f(\bar{x})$ . Equivalently, therefore, we ask that the point  $(\bar{x}, f(\bar{x})) \in \mathbf{R}^n \times \mathbf{R}$  should be close to the affine subspace

$$\{(x, t) \in \mathbf{R}^n \times \mathbf{R} : l_s(x) = t \text{ for all } s \in S\}.$$

As we have observed, the residual in the linear system defining this subspace is small at the point  $(\bar{x}, f(\bar{x}))$ . Standard linear algebra shows that this point is therefore close the solution set (which in particular is nonempty) providing that the smallest singular value of the matrix for the system is bounded below by some fixed tolerance  $\sigma > 0$ . That  $(n+1)$ -by- $k$  matrix has columns  $\begin{pmatrix} \nabla f(s) \\ 1 \end{pmatrix}$  (for  $s \in S$ ), so this lower bound amounts to uniform affine independence of the gradients.

Assuming that this condition holds, success of the algorithm then requires an upper bound on the bundle size  $k$ , since by taking a convergent subsequence of the matrices above we arrive at a limiting matrix with  $k$  linearly independent columns of the form  $\begin{pmatrix} g \\ 1 \end{pmatrix}$ , where each vector  $g$  lies in the limiting gradient set  $\Gamma$  in (3.3), and hence in the subdifferential  $\partial f(\bar{x})$ . Thus the function  $f$  has at least  $k$  affine-independent subgradients at  $\bar{x}$ , from which we deduce the upper bound

$$k \leq 1 + \dim(\partial f(\bar{x})).$$

As we shall see, the uniform affine independence property holds automatically in our convergence proof for functions of max-type. For more general versions of the algorithm, however, we verify the property as follows. For any vectors  $g_i \in \mathbf{R}^n$  indexed by  $i$  in a list  $I$  of length  $k$ , we denote by

$$\sigma_I\{g_i : i \in I\}$$

the  $k$ th largest singular value of a matrix with  $k$  columns  $\begin{pmatrix} g_i \\ 1 \end{pmatrix}$  (for  $i \in I$ ). This nonnegative number is zero exactly when the list is affine dependent. Using this notation, we fix a parameter  $\sigma > 0$  at the outset, and add the following check at the beginning of each iteration.

```

if  $\sigma_S\{\nabla f(s) : s \in S\} < \sigma$  then
  return Stopped: affine dependent gradients;
end if

```

### 3.5 Choosing the bundle size

To summarize, if the algorithm succeeds, the bundle size  $k$  must satisfy both upper and lower bounds involving the Carathéodory number of the limiting gradient set and the dimension of the subdifferential:

$$\text{car } \Gamma \leq k \leq 1 + \dim(\partial f(\bar{x})).$$

In general, these lower and upper bounds may be far apart. For example, for the Euclidean norm  $f = |\cdot|$  at the point  $\bar{x} = 0$ , with the set  $\mathcal{D} = \mathbf{R}^n \setminus \{0\}$ , the bounds become

$$2 \leq k \leq n + 1.$$

However, in the following case of particular interest to us in this work, the two bounds are equal.

**Example 3.5 (Convex max functions)** Consider a nonsmooth function of the form

$$f(x) = \max_{i=1,\dots,k} f_i(x) \quad (x \in \mathbf{R}^n),$$

for smooth convex functions  $f_i: \mathbf{R}^n \rightarrow \mathbf{R}$ , for  $i = 1, 2, \dots, k$ . At the point  $\bar{x} \in \mathbf{R}^n$ , suppose that the function values  $f_i(\bar{x})$  are all equal, with gradients  $g_i = \nabla f_i(\bar{x})$ , so we have

$$\partial f(\bar{x}) = \text{conv}\{g_i : i = 1, 2, \dots, k\}.$$

Assuming that the list of gradients  $\{g_i\}$  is affinely independent, we have

$$k = 1 + \dim(\partial f(\bar{x})).$$

Furthermore, assuming that  $\bar{x}$  is a minimizer, so  $0 \in \partial f(\bar{x})$ , the system

$$\sum_i \lambda_i g_i = 0, \quad \sum_i \lambda_i = 1, \quad \lambda \in \mathbf{R}_+^k$$

must then have a unique solution  $\hat{\lambda} \in \mathbf{R}^k$ . Choose  $\mathcal{D}$  to be the set of points  $x \in \mathbf{R}^n$  for which the maximizing index set  $\text{argmax}_i f_i(x)$  is a singleton, so the limiting gradient set (3.3) is

$$\Gamma = \{g_i : i = 1, 2, \dots, k\}.$$

Hence the Carathéodory number  $\text{car } \Gamma$  is the number of nonzero components of the vector  $\hat{\lambda}$ , which is  $k$  exactly when  $\bar{x}$  is in fact a *nondegenerate* minimizer, meaning that zero lies in the relative interior of the subdifferential  $\partial f(\bar{x})$ .

In general, estimating the lower bound on the bundle size, the Carathéodory number  $\text{car } \Gamma$ , seems challenging in practice. On the other hand, with respect to the upper bound, global nonsmooth optimization methods — the various methods we discussed in the introduction, including the subgradient method, the proximal point and proximal gradient methods, other splitting methods, bundle and level bundle methods, along with nonsmooth BFGS [17], and gradient sampling [3], for example — typically suggest subdifferential dimension information as they progress. Given any finite set of points  $\Omega \subset \mathcal{D}$  near the minimizer  $\bar{x}$ , we can use equation (3.4) to estimate

$$\partial f(\bar{x}) \approx \text{conv}(\nabla f(\Omega)).$$



The dimension of the set on the right-hand side is the rank of a matrix with columns  $(\nabla_1 f(x))$  (for  $x \in \Omega$ ). This suggests that a reasonable estimate of the dimension of  $\partial f(\bar{x})$  is the approximate rank — the number of singular values larger than some tolerance — of this same matrix. This approximate rank then might serve as the integer  $k$  in our Newton method.

### 3.6 The quadratic subproblem

At the end of each iteration we update the bundle  $S$  by substituting a new reference point  $\hat{x} \in \mathbf{R}^n$  for that point in  $S$  whose deletion minimizes the resulting optimality measure. The Newtonian flavor of the algorithm arises from the choice of  $\hat{x} \in \mathbf{R}^n$ , which solves a simple, linearly-constrained quadratic program that we discuss next.

To succeed, the algorithm must generate bundles  $S$  that cluster tightly. By assumption, the corresponding list of gradients  $\nabla f(S)$  is robustly affinely independent, so this gradient information is inconsistent with any smooth model for the objective function  $f$ : we instead must seek a simple, well-behaved, nonsmooth model. We use a max function model, motivated by Example 3.5, and consider twice continuously differentiable functions  $f_s: \mathbf{R}^n \rightarrow \mathbf{R}$  satisfying

$$f_s(s) = f(s), \quad \nabla f_s(s) = \nabla f(s), \quad \nabla^2 f_s(s) = \nabla^2 f(s), \quad f_s(s') < f(s')$$

for distinct points  $s, s' \in S$ . The precise form of these functions is immaterial to the algorithm. It is nonetheless reassuring to note that such functions always exist. We could for example define

$$f_s(x) = q_s(x) - \alpha|x - s|^4 \quad (x \in \mathbf{R}^n),$$

for a sufficiently large constant  $\alpha > 0$ . On the other hand, if in fact  $f$  is a max-function, we could simply consider each  $f_s$  as one of the functions comprising the pointwise maximum.

Now we consider the (unknown) function  $\tilde{f}: \mathbf{R}^n \rightarrow \mathbf{R}$  defined by

$$\tilde{f}(x) = \max_{s \in S} f_s(x) \quad (x \in \mathbf{R}^n),$$

as our working model of the function  $f$ : it agrees with  $f$  up to second order at each of the reference points  $s \in S$ . We can minimize this model via the classical nonlinear program

$$\begin{cases} \text{minimize} & t \\ \text{subject to} & f_s(x) - t \leq 0 \quad (s \in S) \\ & x \in \mathbf{R}^n, \quad t \in \mathbf{R}. \end{cases}$$

Since the functions  $f_s$  are unknown, we cannot solve this problem exactly. Instead, we consider a feasible solution  $(\hat{s}, f(\hat{s}))$ , for any point  $\hat{s} \in \text{conv } S$ , and follow

(loosely) a classical sequential quadratic programming approach to improve it. We remark that SQP techniques have some history in the nonsmooth optimization literature [22].

A standard SQP approach [28] would proceed in two steps, the first of which estimates the Lagrange multipliers. Taking, as a first approximation, each inequality constraint to be active, we seek to solve

$$\min_{\lambda \in \Delta_S} \left| \sum_{s \in S} \lambda_s \nabla f_s(\hat{s}) \right|$$

Approximating the point  $\hat{s}$  by the point  $s$  in each summand leads to the optimality measure in the algorithm:

$$\Theta(S) = \min_{\lambda \in \Delta_S} \left| \sum_{s \in S} \lambda_s \nabla f(s) \right|.$$

Fixing the resulting Lagrange multiplier estimate  $\lambda$ , the Lagrangian for the nonlinear program is

$$(x, t) \mapsto \sum_{s \in S} \lambda_s f_s(x).$$

The second SQP step then aims to reduce its quadratic model at the feasible solution  $(\hat{s}, f(\hat{s}))$  over a feasible region defined by linearized constraints. We approximate the quadratic model by the function

$$(x, t) \mapsto \sum_{s \in S} \lambda_s q_s(x),$$

and we approximate the linearized feasible region, using the active subspace from Section 3.3, as

$$\{(x, t) : l_s(x) = t \text{ for all } s \in S\}.$$

We hence arrive at exactly the quadratic subproblem in the algorithm. The subproblem is feasible, as we saw in Section 3.4, and bounded below by our assumption of strong convexity.

This loose explanation can be tightened. In fact, the proof of Theorem 5.9 will show that, if the point  $\tilde{x} \in \mathbf{R}^n$  minimizes the model  $\tilde{f}$ , and the quantity  $\nu = \max |S - \tilde{x}|$  is small, then, under reasonable conditions the solution  $\hat{x}$  of the quadratic subproblem satisfies  $|\hat{x} - \tilde{x}| = O(\nu^2)$ . In other words, the algorithm computes a good approximation of the minimizer of the model function as the next reference point.

## 4 A sequential quadratic programming tool

For the philosophy underlying the bundle Newton method, the tool we describe in this section is central. It is a slight variant of a standard sequential quadratic

programming technique [32] — for completeness, we prove it directly. Convexity plays no role, throughout this section.

Given functions  $h_i: \mathbf{R}^n \rightarrow \mathbf{R}$ , for  $i = 1, 2, \dots, k$ , we consider an equality-constrained optimization problem of the form

$$(P) \quad \begin{cases} \text{minimize} & c^T y \\ \text{subject to} & h_i(y) = 0 \quad (i = 1, 2, \dots, k) \\ & y \in \mathbf{R}^n. \end{cases}$$

For our purposes, a linear objective function  $c^T y$  (for some vector  $c \in \mathbf{R}^n$ ) suffices.

We say that a point  $\bar{y} \in \mathbf{R}^n$  satisfies the *strong second-order sufficient conditions* if  $h_i(\bar{y}) = 0$  for each  $i$  (feasibility), each  $h_i$  is twice continuously differentiable around  $\bar{y}$ , the list of constraint gradients

$$T = \{\nabla h_i(\bar{y}) : i = 1, 2, \dots, k\}$$

is linearly independent, and there exists a Lagrange multiplier vector  $\bar{\lambda} \in \mathbf{R}^k$  (necessarily unique) satisfying

$$\sum_i \bar{\lambda}_i \nabla h_i(\bar{y}) = -c \quad \text{and} \quad \sum_i \bar{\lambda}_i \nabla^2 h_i(\bar{y}) \text{ positive definite on } T^\perp.$$

The point  $\bar{y}$  is then a strict local minimizer for the problem (P).

We next consider a Lagrange multiplier estimate  $\lambda \in \mathbf{R}^k$  close to  $\bar{\lambda}$ . The traditional approach of sequential quadratic programming linearizes the constraints around a trial point close to the minimizer  $\bar{y}$ , and replaces the objective by the corresponding quadratic approximation to the Lagrangian  $c^T y + \sum_i \lambda_i h_i(y)$ . Instead, we use a *different* reference point  $y_i \in \mathbf{R}^n$  (near  $\bar{y}$ ) for each constraint, forming the corresponding linear approximations  $p_i: \mathbf{R}^n \rightarrow \mathbf{R}$  defined by

$$p_i(y) = h_i(y_i) + \nabla h_i(y_i)^T (y - y_i) \quad (i = 1, 2, \dots, k).$$

We denote by  $Y$  the reference points  $(y_1, y_2, \dots, y_k)$  in the product space  $(\mathbf{R}^n)^k$ , which is close to  $\bar{Y} = (\bar{y}, \bar{y}, \dots, \bar{y})$ , and consider the following quadratic program, parametrized by  $Y$  and  $\lambda$ :

$$(QP) \quad \begin{cases} \text{minimize} & c^T y + \sum_i \lambda_i (p_i(y) + \frac{1}{2}(y - y_i)^T \nabla^2 h_i(y_i)(y - y_i)) \\ \text{subject to} & p_i(y) = 0 \quad (i = 1, 2, \dots, k) \\ & y \in \mathbf{R}^n. \end{cases}$$

The new list of constraint gradients  $\{\nabla h_i(y_i)\}$  is also linearly independent, so any minimizer  $y \in \mathbf{R}^n$  for the quadratic program (QP) must satisfy the conditions for a *stationary point*:

$$\begin{aligned} p_i(y) &= 0 \quad (i = 1, 2, \dots, k) \\ c + \sum_i \lambda_i (\nabla h_i(y_i) + \nabla^2 h_i(y_i)(y - y_i)) &= \sum_i \mu_i \nabla h_i(y_i) \end{aligned}$$

for some multiplier vector  $\mu \in \mathbf{R}^k$ .

**Theorem 4.1** Consider a point  $\bar{y} \in \mathbf{R}^n$  satisfying the strong second-order sufficient conditions for the problem (P). Then for all  $Y \in (\mathbf{R}^n)^k$  near  $\bar{Y}$ , and any multiplier vector  $\lambda = \bar{\lambda} + O(\|Y - \bar{Y}\|)$  in  $\mathbf{R}^k$ , the quadratic program (QP) has a unique stationary point  $\hat{y} = \bar{y} + O(\|Y - \bar{Y}\|^2)$ , which furthermore is the unique minimizer.

**Proof** We can write the stationary point conditions as a linear system:

$$(M(Y, \lambda))(y, \mu) = b(Y, \lambda),$$

for a linear operator  $M(Y, \lambda)$  on  $\mathbf{R}^n \times \mathbf{R}^k$  and a vector  $b(Y, \lambda) \in \mathbf{R}^n \times \mathbf{R}^k$ , both depending continuously on the parameter  $(Y, \lambda)$ . When  $(Y, \lambda) = (\bar{Y}, \bar{\lambda})$ , the corresponding homogeneous system is

$$\begin{aligned} \nabla h_i(\bar{y})^T y &= 0 \quad (i = 1, 2, \dots, k) \\ \sum_i \bar{\lambda}_i \nabla^2 h_i(\bar{y}) y &= \sum_i \mu_i \nabla h_i(\bar{y}). \end{aligned}$$

By the second-order sufficient conditions, this system has only the trivial solution. Hence the operator  $M(\bar{Y}, \bar{\lambda})$  is invertible.

As  $\gamma = \|Y - \bar{Y}\| \rightarrow 0$  with  $|\lambda - \bar{\lambda}| = O(\gamma)$ , we have

$$p_i(\bar{y}) = O(\gamma^2) \quad (i = 1, 2, \dots, k)$$

and

$$\begin{aligned} c + \sum_i \lambda_i (\nabla h_i(y_i) + \nabla^2 h_i(y_i)(\bar{y} - y_i)) &= \sum_i (\lambda_i - \bar{\lambda}_i) \nabla h_i(\bar{y}) + O(\gamma^2) \\ &= \sum_i (\lambda_i - \bar{\lambda}_i) \nabla h_i(y_i) + O(\gamma^2). \end{aligned}$$

We deduce

$$(M(Y, \lambda))(\bar{y}, \lambda - \bar{\lambda}) - b(Y, \lambda) = O(\gamma^2).$$

The norm of the inverse of  $M(Y, \lambda)$  is uniformly bounded for  $(Y, \lambda)$  near  $(\bar{Y}, \bar{\lambda})$ , so

$$(\bar{y}, \lambda - \bar{\lambda}) - (M(Y, \lambda))^{-1}(b(Y, \lambda)) = O(\gamma^2).$$

So there exists a unique stationary point  $\hat{y} = \bar{y} + O(\|Y - \bar{Y}\|^2)$ .

Providing  $\gamma$  is sufficiently small, the point  $\hat{y}$  is in fact a strict minimizer for the quadratic program, since it also satisfies the sufficient condition that  $\sum_i \lambda_i \nabla^2 h_i(y_i)$  is positive definite on the subspace

$$\{z \in \mathbf{R}^n : \nabla h_i(y_i)^T z = 0 \text{ for all } i\}.$$

Otherwise there would exist sequences of points

$$\{y_i^r : i = 1, 2, \dots, k\},$$

and corresponding sequences of multipliers  $\lambda_i^r$  satisfying  $y_i^r \rightarrow \bar{y}$  for each  $i$  and  $\lambda_i^r \rightarrow \bar{\lambda}$  as  $r \rightarrow \infty$ , and unit vectors  $z^r \in \mathbf{R}^n$  satisfying

$$(z^r)^T \left( \sum_i \lambda_i^r \nabla^2 h_i(y_i^r) \right) z^r \leq 0 \quad \text{and} \quad \nabla h_i(y_i^r)^T z^r = 0 \text{ for all } i.$$

In that case, after taking a subsequence, we can suppose that  $z^r$  converges to a unit vector  $z \in \mathbf{R}^n$  satisfying

$$z^T \left( \sum_i \bar{\lambda}_i \nabla^2 h_i(\bar{y}) \right) z \leq 0 \quad \text{and} \quad \nabla h_i(\bar{y})^T z = 0 \text{ for all } i,$$

in contradiction to the second-order sufficient conditions. □

## 5 The max function case

In this section we analyze carefully how the bundle Newton method, Algorithm 2.2, behaves when minimizing a max function  $f: \mathbf{R}^n \rightarrow \mathbf{R}$ , as in Example 3.5. Thus the function  $f$  has the form

$$(5.1) \quad f(x) = \max_{i=1, \dots, k} f_i(x), \quad (x \in \mathbf{R}^n),$$

for some *unknown* functions  $f_i: \mathbf{R}^n \rightarrow \mathbf{R}$  that are now assumed to be twice continuously differentiable for  $i = 1, 2, \dots, k$ .

Like the last section, the development in this section does not rely on convexity, with one key exception. The proof of the final quadratic convergence result, Theorem 5.11 depends critically on strong convexity.

**Note.** We emphasize a crucial point. Our interest in max functions is as local models for more general objectives, and as a test bed for the general-purpose Algorithm 2.2. We could easily minimize an explicitly described objective of the form (5.1) by applying a standard algorithm to the equivalent inequality-constrained optimization problem

$$(IP) \quad \begin{cases} \text{minimize} & t \\ \text{subject to} & f_i(x) - t \leq 0 \quad (i = 1, 2, \dots, k) \\ & x \in \mathbf{R}^n \quad t \in \mathbf{R}. \end{cases}$$

However, we seek to minimize the objective function  $f$  using an oracle *with no access to the individual functions  $f_i$* .

We consider Algorithm 2.2 for the objective  $f$ , on a neighborhood of a point  $\bar{x}$ . Corresponding to the fixed (but implicit) representation (5.1) of  $f$ , we assume the strong second-order conditions defined below.

**Definition 5.2** Given a max function representation of the form (5.1), we say that a point  $\bar{x} \in \mathbf{R}^n$  satisfies the *strong second-order conditions* when the following properties hold.

- *Full activity*: the values  $f_i(\bar{x})$  are equal for all  $i$ .
- *Independence*: the gradients  $\{\nabla f_i(\bar{x}) : i = 1, 2, \dots, k\}$  are affinely independent.
- *Stationarity*: There exists a Lagrange multiplier vector  $\lambda \in \mathbf{R}_+^k$  (necessarily unique) satisfying  $\sum_i \lambda_i = 1$  and  $\sum_i \lambda_i \nabla f_i(\bar{x}) = 0$ .
- *Strict complementarity*:  $\lambda_i > 0$  for all  $i$ .
- *Second-order sufficiency*: The Lagrangian  $\sum_i \lambda_i \nabla^2 f_i(\bar{x})$  is positive definite on the subspace  $\{z \in \mathbf{R}^n : \nabla f_i(\bar{x})^T z \text{ equal for all } i\}$ .

These assumptions are closely related both to problem (IP) and to Example 3.5. If we consider the feasible point  $(\bar{x}, f(\bar{x}))$  for (IP), full activity amounts to all the constraints being active, independence amounts to the usual linear independence constraint qualification, and the remaining three conditions correspond exactly to the analogous conditions for (IP). Classically, these conditions are sufficient for the point  $(\bar{x}, f(\bar{x}))$  to be a strict local minimizer for (IP), and hence for  $\bar{x}$  to be a strict local minimizer for the max function  $f$ . We imposed the first four conditions in Example 3.5: there, we referred to the strict complementarity assumption as non-degeneracy, since, assuming the first three conditions, it amounts to  $0 \in \text{ri}(\partial f(\bar{x}))$  (even in the nonconvex case).

We refer to the disjoint open sets

$$\mathcal{D}_i = \{x \in \mathbf{R}^n : f_i(x) > f_j(x) \ (j \neq i)\},$$

as *activity regions*: the values of the functions  $f$  and  $f_i$  coincide on  $\mathcal{D}_i$ , as do their gradients  $\nabla f$  and  $\nabla f_i$ , and their Hessians  $\nabla^2 f$  and  $\nabla^2 f_i$ . At any point in the open set

$$\mathcal{D} = \bigcup_{i=1}^k \mathcal{D}_i,$$

we suppose that an oracle returns the value of  $f$  along with its gradient and Hessian. Our challenge in minimizing  $f$  is that we only have access to precise information about the value of each component function  $f_i$ , its gradient and Hessian, on the region  $\mathcal{D}_i$ : elsewhere in  $\mathbf{R}^n$  we only know the bound  $f_i \leq f$ .

The algorithm stops if it encounters a point outside  $\mathcal{D}$ . This is a reasonable assumption in practice, since the complement  $\mathcal{D}^c$  is typically a small set. In particular, around the local minimizer  $\bar{x}$ , since the gradients  $\nabla f_i(\bar{x})$  are all distinct (being

affinely independent),  $\mathcal{D}^c$  is contained in the union of the manifolds  $(f_i - f_j)^{-1}(0)$  for  $i \neq j$ . Thus  $\mathcal{D}$  is a dense open set around  $\bar{x}$ .

We consider a closed ball  $B_\gamma(\bar{x})$  of small radius  $\gamma > 0$  around  $\bar{x}$ . At the outset of the algorithm, we consider a *full* bundle  $S \subset B_\gamma(\bar{x})$ , meaning that it contains exactly one reference point in each of the  $k$  activity regions  $\mathcal{D}_i$ . We can therefore write

$$(5.3) \quad S = \{x_1, x_2, \dots, x_k\}, \quad \text{where } x_i \in \mathcal{D}_i \quad i = 1, 2, \dots, k.$$

We will prove that this property is maintained as the algorithm proceeds.

We denote by  $\bar{\sigma}$  a certain  $k$ th largest singular value:

$$(5.4) \quad \bar{\sigma} = \sigma_k \left( \begin{bmatrix} \nabla f_1(\bar{x}) & \nabla f_2(\bar{x}) & \cdots & \nabla f_k(\bar{x}) \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right).$$

The affine independence assumption implies  $\bar{\sigma} > 0$ .

**Proposition 5.5** *For any tolerance  $\sigma \in (0, \bar{\sigma})$ , full bundles  $S$  close to  $\bar{x}$  always satisfy the robust affine independence condition  $\sigma_S(\nabla f(S)) > \sigma$ .*

**Proof** If the result fails, then for each index  $i = 1, 2, \dots, k$  there exists a sequence of points  $(x_i^r)$  in the activity region  $\mathcal{D}_i$  converging to the minimizer  $\bar{x}$ , such that

$$\sigma_k \left( \begin{bmatrix} \nabla f_1(x_1^r) & \nabla f_2(x_2^r) & \cdots & \nabla f_k(x_k^r) \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right) \leq \sigma < \bar{\sigma}$$

for infinitely many  $r$ . But that contradicts the continuity of the  $k$ th largest singular value.  $\square$

We now consider Algorithm 2.2 with the choice of tolerances  $\bar{\epsilon} = 0$ ,  $\bar{\delta} = 0$ , so that the optimality checks never cause the method to stop. (The affine independence check discussed in Section 3.3 will not stop the algorithm either, providing we fix the tolerance  $\sigma \in (0, \bar{\sigma})$ .) We begin each iteration by forming the corresponding linear and quadratic approximations:

$$\begin{aligned} l_i(\cdot) &= f_i(x_i) + \nabla f_i(x_i)^T(\cdot - x_i), \\ q_i(\cdot) &= l_i(\cdot) + \frac{1}{2}(\cdot - x_i)^T \nabla^2 f_i(x_i)(\cdot - x_i), \end{aligned}$$

and estimate the Lagrange multiplier vector. We have the following result.

**Proposition 5.6** *For any small radius  $\gamma > 0$ , there exists a unique minimizer for the problem*

$$\min_{\lambda \in \mathbf{R}^k} \left\{ \left| \sum_i \lambda_i \nabla f_i(x_i) \right| : \sum_i \lambda_i = 1 \right\},$$

and it satisfies  $\lambda = \bar{\lambda} + O(\gamma)$ .

**Proof** A vector  $\lambda \in \mathbf{R}^k$  solves the problem if and only if there exists a number  $\alpha \in \mathbf{R}$  such that

$$\sum_{i=1}^k \lambda_i = 1 \quad \text{and} \quad \alpha + \nabla f_j(x_j)^T \left( \sum_{i=1}^k \lambda_i \nabla f_i(x_i) \right) = 0 \quad (j = 1, 2, \dots, k).$$

This square linear system is defined by an operator depending smoothly on its parameters, the points  $x_i$  (for  $i = 1, 2, \dots, k$ ). In the limit, when  $x_i = \bar{x}$  for each  $i$ , the system becomes

$$\sum_{i=1}^k \lambda_i = 1 \quad \text{and} \quad \alpha + \nabla f_j(\bar{x})^T \left( \sum_{i=1}^k \lambda_i \nabla f_i(\bar{x}) \right) = 0 \quad (j = 1, 2, \dots, k).$$

The corresponding homogeneous system has only the trivial solution, by affine independence of the set  $\{\nabla f_i(\bar{x})\}$ , so the defining operator is invertible. The unique solution of this limiting system is clearly  $(\bar{\lambda}, 0)$ , and the result now follows.  $\square$

As a consequence of this result and the strict complementarity assumption, the Lagrange multiplier estimates are all positive for small radius  $\gamma$ . We next turn to the computation of the new reference point  $\hat{x}$ .

Being a fully active strict local minimizer of the inequality-constrained problem (IP), the point  $(\bar{x}, f(\bar{x}))$  is also a local minimizer for the more restrictive problem

$$(P') \quad \begin{cases} \text{minimize} & t \\ \text{subject to} & f_i(x) - t = 0 \quad (i = 1, 2, \dots, k) \\ & x \in \mathbf{R}^n \quad t \in \mathbf{R}. \end{cases}$$

Furthermore, as is easy to verify, it satisfies the strong second-order sufficient conditions with the same Lagrange multiplier vector  $\bar{\lambda}$ . Hence we can apply Theorem 4.1. The corresponding quadratic program is

$$\begin{cases} \text{minimize} & t + \sum_i \lambda_i (q_i(x) - t) \\ \text{subject to} & l_i(x) - t = 0 \quad (i = 1, 2, \dots, k) \\ & x \in \mathbf{R}^n, t \in \mathbf{R}, \end{cases}$$

or equivalently, exactly the quadratic program in Algorithm 2.2:

$$\begin{cases} \text{minimize} & \sum_i \lambda_i q_i(x) \\ \text{subject to} & l_i(x) \text{ equal for all } i = 1, 2, \dots, k \\ & x \in \mathbf{R}^n. \end{cases}$$

By Theorem 4.1, this quadratic program has a unique minimizer  $\hat{x} \in \mathbf{R}^n$ , which furthermore satisfies  $\hat{x} = \bar{x} + O(\gamma^2)$ .

The final step in the iteration substitutes the new reference point  $\hat{x}$  for one of the existing reference points. Since  $\hat{x}$  is also close to the point  $\bar{x}$ , this substitution produces another full bundle. To see this, the following idea is the key tool.



**Proposition 5.7** *Near  $\bar{x}$ , the optimality measures of full bundles are always strictly less than those of bundles that are not full.*

**Proof** If the result fails, then there exists a sequence of full bundles  $S_r$ , and a sequence of not full bundles  $S'_r$ , both shrinking to  $\bar{x}$ , with

$$\Theta(S_r) \geq \Theta(S'_r) \quad (r = 1, 2, 3, \dots).$$

The left-hand side converges to zero, because

$$0 \in \text{conv}\{\nabla f_i(\bar{x}) : i = 1, 2, \dots, k\},$$

and hence so does the right-hand side. After taking a subsequence, we can suppose that there is an index  $j$  such that  $S'_r \cap \mathcal{D}_j = \emptyset$ , and hence

$$\liminf_r \Omega(S'_r) \geq \min |\text{conv}\{\nabla f_i(\bar{x}) : i \neq j\}|.$$

But the right-hand side is strictly positive, because

$$0 \notin \text{conv}\{\nabla f_i(\bar{x}) : i \neq j\}.$$

This contradiction completes the proof.  $\square$

Since the reference points stay close to the point  $\bar{x}$ , we next deduce that we maintain full bundles as the algorithm progresses, as follows.

**Corollary 5.8** *For any full bundle  $S$  near  $\bar{x}$ , and any new reference point  $\hat{x} \in \mathcal{D}$  near  $\bar{x}$ , there is a unique reference point  $s \in S$  minimizing the optimality measure of the new bundle  $S' = (S \setminus \{s\}) \cup \{\hat{x}\}$ , and  $S'$  is then also a full bundle.*

**Proof** The previous result shows that when  $\hat{x}$  lies in the activity region  $\mathcal{D}_i$ , the unique optimal choice of  $s$  is the unique reference point in  $\mathcal{D}_i$ . The result then follows.  $\square$

In summary, we have proved the following result.

**Theorem 5.9** *Given a max function representation of the objective*

$$f(x) = \max_{i=1, \dots, k} f_i(x) \quad (x \in \mathbf{R}^n),$$

*suppose the point  $\bar{x}$  satisfies the strong second-order conditions in Definition 5.2. Then there exists a constant  $M > 0$  such that the  $k$ -bundle Newton method (Algorithm 2.2), with the tolerances  $\bar{\epsilon} = 0$ ,  $\bar{\delta} = 0$ , has the following property. Any iteration starting with a full bundle  $S$  sufficiently close to  $\bar{x}$  generates a new reference point  $\hat{x}$  satisfying*

$$|\hat{x} - \bar{x}| \leq M \max_{s \in S} |s - \bar{x}|^2,$$

*and assuming  $\hat{x} \in \mathcal{D}$ , then generates a new full bundle by substituting  $\hat{x}$  for the unique reference point in  $S$  from the same activity region.*

While this is a suggestive result, it does not yet guarantee convergence. To ensure that the sequence of bundles shrinks to  $\bar{x}$ , we finally call on our assumption that the objective  $f$  is strongly convex. (In the next section we discuss how to modify the algorithm to handle nonconvex objectives.)

We first develop a simple tool. As usual, for vectors  $z \in \mathbf{R}^k$  we define  $\|z\|_{\max} = \max_j |z_j|$ .

**Lemma 5.10** *Given constants  $\epsilon, M > 0$ , consider any sequence of vectors in the orthant  $\mathbf{R}_+^n$  with the property that, for each successive pair  $z, z'$  in the sequence, there exists an index  $i$  such that  $z'_j = z_j$  for all  $j \neq i$ , and furthermore*

$$z_i \geq \epsilon \|z\|_{\max} \quad \text{and} \quad z'_i \leq M \|z\|_{\max}^2.$$

*Then providing the initial vector is sufficiently small, the sequence converges to zero at a  $k$ -step quadratic rate.*

**Proof** By induction we see that  $\|z\|_{\max}$  is nondecreasing as the vector  $z$  evolves along the sequence, providing that the initial vector is sufficiently small. Suppose  $z = z_{\text{old}}$  at the outset of some iteration, and set  $\theta = \|z_{\text{old}}\|_{\max}$ . At this and every subsequent iteration, the updated component  $z_i$  is always set to a new value in the interval  $(0, M\theta^2]$ . Each updated component therefore cannot be updated again until we have

$$\max_{j=1,2,\dots,k} z_j \leq \frac{M\theta^2}{\epsilon}.$$

This inequality must therefore hold after at most  $k$  iterations, at which point, if  $z = z_{\text{new}}$ , we have

$$\|z_{\text{new}}\|_{\max} \leq \frac{M}{\epsilon} \|z_{\text{old}}\|_{\max}^2,$$

which completes the proof. □

We can now prove our main result.

**Theorem 5.11 (Fast convergence for strongly convex max functions)**

*Given a max function representation of the objective*

$$f(x) = \max_{i=1,\dots,k} f_i(x) \quad (x \in \mathbf{R}^n),$$

*suppose that the point  $\bar{x}$  satisfies the strong second-order conditions in Definition 5.2, with each Hessian  $\nabla^2 f_i(\bar{x})$  positive definite. Then, given the tolerances  $\bar{\epsilon} = 0$ ,  $\bar{\delta} = 0$ , the  $k$ -bundle Newton method (Algorithm 2.2) starting from any full bundle in a neighborhood of  $\bar{x}$ , either stops at a point outside the set  $\mathcal{D}$ , or generates a sequence of full bundles that converge  $k$ -step quadratically to  $\bar{x}$ .*

**Proof** Given any small radius  $\gamma > 0$ , each function  $f_i$  is  $\rho$ -strongly convex on  $B_\gamma(\bar{x})$ , for some constant  $\rho > 0$ . In fact, the proof that follows takes place entirely in the ball  $B_\gamma(\bar{x})$ , so we lose no generality in assuming that each function  $f_i$  is convex.

According to Theorem 5.9, there exists a constant  $M > 0$  such that, for small enough  $\gamma > 0$ , starting from any full bundle in the ball  $B_\gamma(\bar{x})$ , Algorithm 2.2 produces a sequence of full bundles

$$S = \{x_1, x_2, \dots, x_k\}, \quad \text{where } x_i \in \mathcal{D}_i \cap B_\gamma(\bar{x}) \quad i = 1, 2, \dots, k,$$

and at each iteration replaces a reference point,  $x_i$  for some index  $i$  with a new reference point  $\hat{x} \in \mathcal{D}_i \cap B_\gamma(\bar{x})$  which furthermore satisfies

$$(5.12) \quad |\hat{x} - \bar{x}| \leq M \max_{j=1, \dots, k} |x_j - \bar{x}|^2.$$

By the construction, the new reference point  $\hat{x}$  satisfies

$$l_j(\hat{x}) = l_i(\hat{x}) \quad \text{for } j = 1, 2, \dots, k,$$

and since the functions  $f_i$  are twice continuously differentiable, there exists a constant  $R > 0$  such that

$$f_i(x) \leq l_i(x) + \frac{R}{2}|x - x_i|^2 \quad \text{for all } x \in B_\gamma(\bar{x}).$$

On the other hand, by strong convexity we have

$$f_j(x) \geq l_j(x) + \frac{\rho}{2}|x - x_j|^2 \quad \text{for all } x \in B_\gamma(\bar{x}), \quad j = 1, 2, \dots, k.$$

Therefore we deduce

$$l_i(\hat{x}) + \frac{R}{2}|\hat{x} - x_i|^2 \geq f_i(\hat{x}) = f_j(\hat{x}) \geq l_j(\hat{x}) + \frac{\rho}{2}|\hat{x} - x_j|^2 = l_i(\hat{x}) + \frac{\rho}{2}|\hat{x} - x_j|^2,$$

and setting  $\alpha = \sqrt{\rho/R} > 0$  gives

$$|\hat{x} - x_i| \geq \alpha|\hat{x} - x_j| \quad \text{for } j = 1, 2, \dots, k.$$

Let  $\beta = \max_j |x_j - \bar{x}|$ , so by inequality (5.12) we have  $|\hat{x} - \bar{x}| \leq M\beta^2$ . The preceding inequality implies, for each  $j = 1, 2, \dots, k$ , the inequality

$$\begin{aligned} \alpha|x_j - \bar{x}| &\leq \alpha|x_j - \hat{x}| + \alpha|\hat{x} - \bar{x}| \leq |x_i - \hat{x}| + M\alpha\beta^2 \\ &\leq |x_i - \bar{x}| + |\bar{x} - \hat{x}| + M\alpha\beta^2 \leq |x_i - \bar{x}| + M(\alpha + 1)\beta^2 \end{aligned}$$

Maximizing over  $j$  implies

$$|x_i - \bar{x}| \geq \alpha\beta - M(\alpha + 1)\beta^2 \geq \frac{\alpha}{2}\beta$$

providing the radius  $\gamma$  is sufficiently small. So we conclude that the old reference point  $x_i$  and the new reference point  $\hat{x}$  satisfy the two key inequalities

$$|x_i - \bar{x}| \geq \frac{\alpha}{2} \max_j |x_j - \bar{x}| \quad \text{and} \quad |\hat{x} - \bar{x}| \leq M \max_j |x_j - \bar{x}|^2.$$

We now define

$$z_j = |x_j - \bar{x}| \in (0, \gamma) \quad (j = 1, 2, \dots, k),$$

and apply Lemma 5.10 to complete the proof.  $\square$

Notice that the final assumption in the strong second-order conditions — that the Hessian of the Lagrangian is positive definite on the tangent subspace — is in fact superfluous here, since we are assuming that each Hessian  $\nabla^2 f_i(\bar{x})$  is positive definite.

## 6 Minimizing smooth-nonsmooth sums

With easy modifications, we can use the same Newton bundle idea as in Algorithm 2.2 to minimize a function of the form  $F = f + r$ , where the first component  $f: \mathbf{R}^n \rightarrow \mathbf{R}$  is strongly convex but nonsmooth as before, and the second component  $r: \mathbf{R}^n \rightarrow \mathbf{R}$  is nonconvex but smooth (twice continuously differentiable).

Using the natural choice of optimality measure for a bundle  $S \subset \mathcal{D}$ , namely

$$(6.1) \quad \Theta(S) = \min |\text{conv}(\nabla F(S))|.$$

the new algorithm proceeds as follows.

**Algorithm 6.2** (*k*-bundle Newton method for  $F = f + r$ )

**Require:** initial bundle  $S \subset \mathcal{D}$  of size  $k$ , tolerances  $\bar{\epsilon}, \bar{\delta} \geq 0$ ;

**for** iteration = 1, 2, 3, ... **do**

**for**  $s \in S$  **do**

$$l_s(\cdot) = f(s) + \nabla f(s)^T(\cdot - s);$$

$$q_s(\cdot) = F(s) + \nabla F(s)^T(\cdot - s) + \frac{1}{2}(\cdot - s)^T \nabla^2 F(s)(\cdot - s);$$

**end for**

$$\delta = \Theta(S);$$

  choose  $\lambda \in \Delta_S$  with  $|\sum_{s \in S} \lambda_s \nabla F(s)| = \delta$ ;

**if**  $\text{diam } S < \bar{\epsilon}$  and  $\delta < \bar{\delta}$  **then**

**return** *Stopped: nearly optimal*;

**end if**

  choose  $\hat{x} \in \text{argmin} \{ \sum_{s \in S} \lambda_s q_s(x) : x \in \mathbf{R}^n, l_s(x) \text{ equal for all } s \in S \}$ ;

**if**  $\hat{x} \notin \mathcal{D}$  **then**

**return** *Stopped: nonsmooth point.*;

```

else
  choose  $s \in S$  minimizing  $\Theta((S \setminus \{s\}) \cup \{\hat{x}\})$ ;
   $S = (S \setminus \{s\}) \cup \{\hat{x}\}$ ;
end if
end for

```

We make some immediate comments in comparison with Algorithm 2.2 for minimizing  $f$  alone. Important to notice is that the affine functions  $l_s$  (for  $s \in S$ ), whose equality defines the subproblem constraints, are linear approximations for the component  $f$ , rather than for the objective  $F$ : this subtlety is important for the analysis, which relies on strong convexity of  $f$ . In contrast, the functions  $q_s$  appearing in the subproblem objective must be quadratic approximations for  $F$  to ensure quadratic convergence.

In this new setting we should consider the possibility of an unbounded quadratic subproblem, since the quadratic approximations are no longer necessarily convex. For the time being, to emphasize the parallel with Algorithm 2.2, we omit this check, but include it in the final version (still to come).

As in Algorithm 2.2, we might also check affine independence of the gradients: whether we use the gradients of the component  $f$  or the objective  $F$  is immaterial, since for bundles of small diameter the two differ by roughly a constant due to the continuity of the gradient of the other component  $r$ . When  $f$  is a max function, the analogue of Proposition 5.5 still holds: providing we choose a sufficiently small affine dependence parameter, full bundles close to  $\bar{x}$  always satisfy the check. Like the unboundedness check, we omit this check for now, but include it in the final version.

The optimality measure (6.1) is the natural extension of the earlier version, and leads to the precise analogue of the approximate optimality property (3.2), as well as to our strategy for substituting the new reference point at the end of each iteration.

All of our discussion of the choice of bundle size applies in this new case too, simply replacing the old objective  $f$  by the new objective  $F$ , and noting the simple relationship between the subdifferentials

$$\partial F(x) = \partial f(x) + \nabla r(x) \quad (x \in \mathbf{R}^n),$$

where we understand the left-hand side in the Clarke sense [4]. To motivate the quadratic subproblem, we follow exactly the same argument as in Section 3.6, but modifying the underlying nonlinear program to

$$\begin{cases} \text{minimize} & r(x) + t \\ \text{subject to} & f_s(x) - t \leq 0 \\ & x \in \mathbf{R}^n, \quad t \in \mathbf{R}. \end{cases} \quad (s \in S)$$

We note that the constraints are not changed: they do not involve the function  $r$ . Hence, as we noted, neither does the active subspace in the algorithm.

We can then mimic the analysis of the original bundle-Newton method for max functions. For notational simplicity, we proved the sequential quadratic programming tool, Theorem 4.1, for a linear objective, but the case we now need, involving a nonlinear objective, is almost identical.

We thus arrive at the main analysis, Section 5. Given a representation of the nonsmooth component function  $f$  in the objective as

$$f(x) = \max_{i=1,\dots,k} f_i(x) \quad (x \in \mathbf{R}^n),$$

we assume that the point  $\bar{x} \in \mathbf{R}^n$  satisfies the *strong second-order conditions* for the problem of minimizing the objective  $f+r$ . By this, we mean that the feasible point  $(\bar{x}, f(\bar{x}))$  for the appropriate modification of the problem (IP), namely,

$$\begin{cases} \text{minimize} & r(x) + t \\ \text{subject to} & f_i(x) - t \leq 0 \quad (i = 1, 2, \dots, k) \\ & x \in \mathbf{R}^n \quad t \in \mathbf{R}, \end{cases}$$

has all constraints active, and satisfies the linear independence constraint qualification, stationarity, strict complementarity, and the second-order sufficient conditions. Once again we note that the constraints are unchanged: they do not involve the function  $r$ . With a virtually identical proof, we arrive at the following generalization of Theorem 5.11.

**Theorem 6.3 (Quadratic convergence for smooth-nonsmooth sums)**

*Given a max function representation of the nonsmooth function*

$$f(x) = \max_{i=1,\dots,k} f_i(x) \quad (x \in \mathbf{R}^n),$$

*and a twice continuously differentiable function  $r: \mathbf{R}^n \rightarrow \mathbf{R}$ , suppose that the point  $\bar{x} \in \mathbf{R}^n$  satisfies the strong second-order conditions for minimizing the objective  $f+r$ , with each Hessian  $\nabla^2 f_i(\bar{x})$  positive definite. Then, given the tolerances  $\bar{\epsilon} = 0$  and  $\bar{\delta} = 0$ , the  $k$ -bundle Newton method (Algorithm 6.2) starting from any full bundle in a neighborhood of  $\bar{x}$ , either stops at a point outside the set  $\mathcal{D}$ , or generates a sequence of full bundles that converge  $k$ -step quadratically to  $\bar{x}$ .*

## 7 Minimizing weakly convex functions

Consider a *weakly convex* function  $F: \mathbf{R}^n \rightarrow \mathbf{R}$ , by which we mean that, for some *weak convexity parameter* value  $\eta$ , the function  $f = F + \frac{\eta}{2}|\cdot|^2$  is convex. By increasing  $\eta$  if necessary,  $f$  becomes strongly convex. Assuming that  $F$  is also twice continuously differentiable around every point in the set  $\mathcal{D}$ , we can then define the smooth function  $r = -\frac{\eta}{2}|\cdot|^2$  and arrive at the following version of Algorithm 6.2 for minimizing  $F$ , using the optimality measure:

$$\Theta(S) = \min |\text{conv}(\nabla F(S))|.$$

**Algorithm 7.1** (*k*-bundle Newton minimization for weakly convex  $F$ )

**Require:** initial bundle  $S \subset \mathcal{D}$  of size  $k$ , tolerances  $\bar{\epsilon}, \bar{\delta} \geq 0$ ,  $\sigma > 0$ ,

weak convexity parameter  $\eta$ ;

**for** iteration = 1, 2, 3, ... **do**

**for**  $s \in S$  **do**

$$l_s(\cdot) = F(s) + \frac{\eta}{2}|s|^2 + (\nabla F(s) + \eta s)^T(\cdot - s);$$

$$q_s(\cdot) = F(s) + \nabla F(s)^T(\cdot - s) + \frac{1}{2}(\cdot - s)^T \nabla^2 F(s)(\cdot - s);$$

**end for**

**if**  $\sigma_S\{\nabla F(s) : s \in S\} < \sigma$  **then**

**return** *Stopped: affine dependent gradients.*

**end if**

$\delta = \Theta(S)$ ;

  choose  $\lambda \in \Delta_S$  with  $|\sum_{s \in S} \lambda_s \nabla F(s)| = \delta$ ;

**if**  $\text{diam } S < \bar{\epsilon}$  and  $\delta < \bar{\delta}$  **then**

**return** *Stopped: nearly optimal;*

**end if**

**if**  $\min\{\sum_{s \in S} \lambda_s q_s(x) : x \in \mathbf{R}^n, l_s(x) \text{ equal for all } s \in S\} = -\infty$  **then**

**return** *Stopped: unbounded subproblem.*

**else**

    choose optimal  $\hat{x}$ ;

**if**  $\hat{x} \notin \mathcal{D}$  **then**

**return** *Stopped: nonsmooth point.;*

**else**

      choose  $s \in S$  minimizing  $\Theta((S \setminus \{s\}) \cup \{\hat{x}\})$ ;

$$S = (S \setminus \{s\}) \cup \{\hat{x}\};$$

**end if**

**end if**

**end for**

We note that the algorithm is almost identical to the original convex version, Algorithm 2.2, the only change (other than termination checks) being to the definition of the linear approximations defining the active subspace. As the weak convexity parameter  $\eta$  grows large, the active subspace converges to the subspace of all points equidistant from all the reference points in the bundle.

Theorem 6.3 then specializes to our culminating result. Recall that the constant  $\bar{\sigma} > 0$  is defined by equation (5.4).

**Corollary 7.2** (Fast convergence for weakly convex max functions)

*Given a max function representation of the objective*

$$F(x) = \max_{i=1, \dots, k} f_i(x) \quad (x \in \mathbf{R}^n),$$

*suppose that the point  $\bar{x} \in \mathbf{R}^n$  satisfies the strong second-order conditions given in Definition 5.2. Then, given the tolerances  $\bar{\epsilon} = 0$ ,  $\bar{\delta} = 0$ , small  $\sigma > 0$ , and sufficiently*

large weak convexity parameter  $\eta$ , the  $k$ -bundle Newton method (Algorithm 7.1) starting from any full bundle in a neighborhood of  $\bar{x}$ , either stops at a point outside the set  $\mathcal{D}$ , or generates a sequence of full bundles that converge  $k$ -step quadratically to  $\bar{x}$ .

**Note:** Any choice of parameters  $\sigma \in (0, \bar{\sigma})$  and  $\eta$  strictly larger than the largest eigenvalues of each negative Hessian  $-\nabla^2 f_i(\bar{x})$  in fact suffices.

**Proof** This follows directly from Theorem 6.3, once we observe that, since each function  $f_i$  is twice continuously differentiable, with the given choice of weak convexity parameter  $\eta$  the functions  $f_i + \frac{\eta}{2}|\cdot|^2$  are all strongly convex on a neighborhood of  $\bar{x}$ , and this local property suffices for the proof.  $\square$

## 8 Numerical experiments

We illustrate the local bundle Newton method on several nonsmooth objective functions. This handful of simple experiments is meant as a proof of concept rather than comprising any algorithmic recommendations. Nonetheless, the results appear clearly promising enough to invite future research.

### 8.1 Practical considerations

We implemented none of the stopping criteria, simply terminating the algorithm manually when rounding error prevented any further progress.

#### Choosing an initial bundle

In each experiment, we ran a standard global nonsmooth optimization method to generate a finite set of points  $\Omega \subset \mathcal{D}$  near a minimizer  $\bar{x}$ , and used the corresponding gradients to estimate the dimension of the subdifferential  $\partial f(\bar{x})$  and hence choose the bundle size  $k$ , as we discussed in Section 3.5. We then ran a heuristic subset selection procedure [8] to choose a set of  $k$  points in  $\Omega$  with robustly affinely independent gradients to form the initial bundle.

For convex problems, we implemented the simple “Bundle Method with Multiple Cuts” [6], which we reproduce below in our notation. In the nonconvex case, we implemented the nonsmooth BFGS method [17].

#### Algorithm 8.1 (Multiple cut bundle method to minimize convex $f$ )

**Require:** initial bundle  $S \subset \mathbf{R}^n$ , initial center  $z \in S$ , stopping tolerance  $\bar{\epsilon}$ , proximal parameter  $\rho > 0$ , sufficient decrease parameter  $\beta \in (0, 1)$   
**for** iteration = 1, 2, 3, ... **do**



```

for  $s \in S$  do
   $g_s \in \partial f(s)$ ;
   $l_s(\cdot) = f(s) + g_s^T(\cdot - s)$ ;
end for
Choose  $\hat{x}$  minimizing  $\max_{s \in S} l_s(\cdot) + \frac{\rho}{2} |\cdot - z|^2$ ;
if  $f(z) - \max_{s \in S} l_s(\hat{x}) \leq \bar{\epsilon}$  then
  return Stopped: nearly optimal.
else
  if  $f(\hat{x}) \leq f(z) - \beta(f(z) - \max_{s \in S} l_s(\hat{x}))$  then
     $z \leftarrow \hat{x}$  (serious step)
  else
     $z \leftarrow z$  (null step)
  end if
   $S \leftarrow S \cup \{\hat{x}\}$ ;
end if
end for

```

For the bundle method, we chose  $\Omega$  to be the set of points whose cutting planes were strongly active in the final iteration. That is,

$$\Omega = \{s \in S : \alpha_s > 0\},$$

where  $\alpha_s$  is the dual variable associated with cutting plane  $l_s(\cdot)$ . For BFGS, we chose the set  $\Omega$  to be the final  $2n$  iterates.

### Solving the quadratic subproblems

The algorithm involves two quadratic programming subproblems. The first involves computing the optimality measure  $\Theta(S)$  which amounts to projecting 0 onto the convex hull of vectors  $\{\nabla f(s) : s \in S\}$ . We implemented this as a quadratic program, solved in Gurobi. For the equality-constrained quadratic programs,

$$(8.2) \quad \begin{cases} \text{minimize} & \sum_s \lambda_s q_s(x) \\ \text{subject to} & l_s(x) - t = 0 \quad (s \in S) \\ & t \in \mathbf{R}, \quad x \in \mathbf{R}^n. \end{cases}$$

we simply solve the (linear) optimality conditions directly:

$$(8.3) \quad \begin{aligned} \sum_{s \in S} \lambda_s \nabla^2 f(s)(x - s) + \sum_{s \in S} \mu_s \nabla f(s) &= 0, \\ \sum_{s \in S} \mu_s &= 1, \\ l_s(x) - t &= 0 \quad (s \in S). \end{aligned}$$

The  $x$  variable of the solution is then our Newton iterate  $\hat{x}$ .

## 8.2 Illustrative Examples

### A Strongly Convex Problem

Our first experiment is to minimize max functions of the form

$$(8.4) \quad f(x) = \max_{i=1,\dots,k} \left\{ g_i^T x + \frac{1}{2} x^T H_i x + \frac{c_i}{24} \|x\|^4 \right\}$$

for  $1 \leq k \leq n+1$ . We randomly generate positive constants  $c_i$ , symmetric positive definite matrices  $H_i$ , and affinely independent random vectors  $g_i$  satisfying  $\sum_i \lambda_i g_i = 0$  for some  $\lambda$  randomly sampled in  $\{\lambda > 0 : \sum_i \lambda_i = 1\}$ . Then

$$0 \in \partial f(0) = \text{conv}\{g_i : 1 \leq i \leq k\}$$

so  $f$  is nonsmooth at the minimizer of 0. The structure of  $f$  is unknown to the algorithms, whose access is limited to a black box that returns function values, gradients, and Hessians.

In random trials for dimension  $n = 50$ , we applied the bundle method, Algorithm 8.1, in a first phase, with parameters  $\rho = 1$  and  $\beta = 10^{-5}$  and starting point  $z = (1, \dots, 1)$ . The stopping tolerance was set to  $10^{-6}$ , at which point we initialized and switched to the bundle Newton method, Algorithm 2.2. Results for a number of random trials are shown in Figures 1 and 2. For the bundle method phase, we observed a roughly linear rate of convergence of function values to zero, proportional to the degree of nonsmoothness  $k$ . Switching to the bundle Newton method results in much more rapid convergence in accordance with the rates predicted by Theorem 5.11.

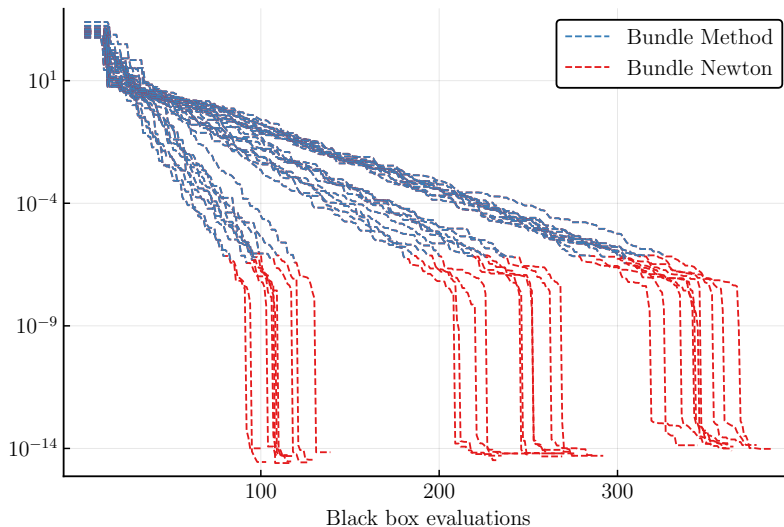


Figure 1: Best function value found for the bundle method and bundle Newton method against number of black box calls for random max functions (8.4) for  $k = 10, 25, 40$  in dimension  $n = 50$ .

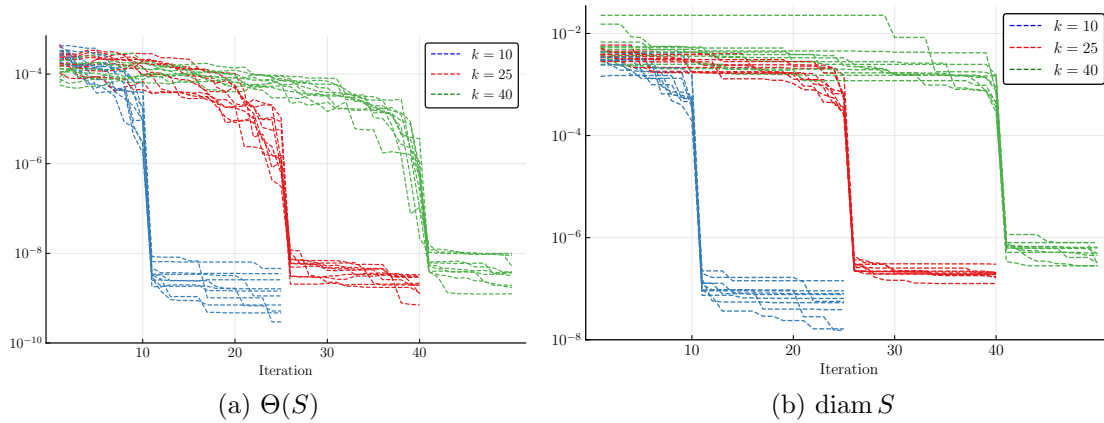


Figure 2: Optimality measures against iteration count for bundle Newton method for random max functions (8.4).

### A Nonconvex Problem

To test the nonconvex version of the algorithm, we used Euclidean sum functions of the form

$$(8.5) \quad f(x) = \sum_{i=1}^k \left| g_i^T x + \frac{1}{2} x^T H_i x + \frac{c_i}{24} \|x\|^4 \right|$$

for  $1 \leq k \leq n + 1$ . The constants  $c_i$ , vectors  $g_i$  and matrices  $H_i$  were randomly generated as in the previous experiment. As usual, access to  $f$  was limited to a black box that returns function values, gradients, and Hessians.

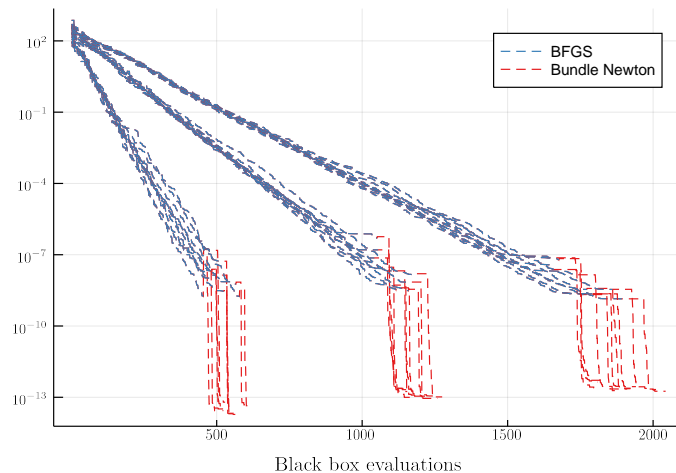


Figure 3: Best function value found for BFGS and the bundle Newton method against number of black box calls on random Euclidean sum functions (8.5) for  $k = 10, 25, 40$  in dimension  $n = 50$ .

In random trials for dimension  $n = 50$ , we applied nonsmooth BFGS in a first phase until a breakdown occurred due to numerical instability (as usual with this method [17]). At this point we switched to the Algorithm 7.1 with weak convexity parameter dynamically chosen as

$$\eta = \max_{s \in S} \lambda_{\max}(-\nabla^2 f(s))$$

at each iteration.

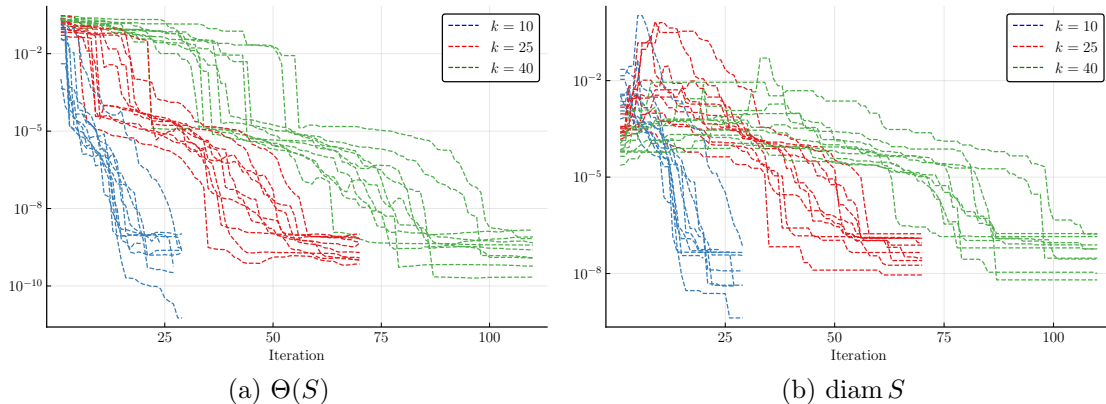


Figure 4: Optimality measures against iteration count for bundle Newton method for random Euclidean sum functions (8.5).

### 8.3 Partly Smooth Functions

While the theoretical results of this paper were limited to objective functions with finite max structure, experiments suggest that variants of the bundle Newton method may be effective much more broadly. However, one particular implementation hurdle arises even for simple nonsmooth functions like the Euclidean norm: solving the system (8.3) directly will be numerically unstable, due to the ill-conditioning of the Hessians  $\nabla^2 f(s)$ .

However, for the broad class of *partly smooth* functions in the sense of [16], this ill-conditioning is often highly structured: the nonsmoothness is associated with a certain subspace  $\mathcal{V}$  spanned by the eigenvectors corresponding to large eigenvalues of the Hessian  $\nabla^2 f(s)$ . Along an orthogonal manifold, the function behaves smoothly, with well conditioned Hessian. Motivated by this idea, we follow a simple strategy for solving the system (8.3), similar to reduced system approaches for nonlinear programming described in standard texts [28], and avoiding full Hessian computations.

## A reduced system approach

Let  $G$  and  $b$  be a matrix and vector satisfying

$$\{x : Gx = b\} = \{x : l_s(x) \text{ equal for all } s \in S\},$$

so that we can write the optimality conditions (8.3) as

$$(8.6) \quad \begin{aligned} \sum_{s \in S} \lambda_s \nabla^2 f(s)(x - s) + G^T \nu &= - \sum_{s \in S} \lambda_s \nabla f(s), \\ Gx &= b. \end{aligned}$$

for  $\nu \in \mathbf{R}^{k-1}$ . Suppose that we have found matrices  $U$  and  $V$  such that the matrix  $\begin{bmatrix} U & V \end{bmatrix} \in \mathbf{R}^{n \times n}$  is full rank and  $GU = 0$  (via a QR factorization of  $G^T$ , for example). The columns of  $U$  are then a basis for the space  $\text{Null}(G)$ , and we can write any solution of (8.6) as  $x = Ux_u + Vx_v$ . The constraint  $Gx = b$  then implies  $GVx_v = b$ , which can be solved for  $x_v$ , since we assume  $G$  (and hence  $GV$ ) is full rank. We deduce

$$\{x : Gx = b\} = \text{Range}(U) + p,$$

where  $p$  is the particular solution  $V(GV)^{-1}b$ . Substituting this into the stationarity condition and multiplying through by  $U^T$  yields the *reduced* system

$$\sum_{s \in S} \lambda_s U^T \nabla^2 f(s)(Ux_u + p - s) = - \sum_{s \in S} \lambda_s U^T \nabla f(s).$$

In a slight modification to the algorithm, if we project each reference point onto the active subspace we arrive at the linear system

$$\sum_{s \in S} \lambda_s U^T \nabla^2 f(s) U x_u = \sum_{s \in S} \lambda_s [(U^T \nabla^2 f(s) U) U^T (s - p) - U^T \nabla f(s)].$$

This system only involves the *projected Hessians*  $U^T \nabla^2 f(s) U$ , which remain well conditioned if the span of  $V$  is close to the subspace  $\mathcal{V}$ , a property that we have experimentally observed to hold in practice.

## An Eigenvalue Problem

Our final experiment to illustrate the reduced systems approach is an eigenvalue problem. Specifically, given symmetric matrices  $A_0, \dots, A_n \in \mathbf{R}^{m \times m}$  we seek to minimize

$$(8.7) \quad f(x) = \lambda_{\max} \left( A_0 + \sum_{i=1}^n x_i A_i \right),$$

where  $\lambda_{\max}(\cdot)$  is the largest eigenvalue function. Typically minimizers occur at points where  $\lambda_{\max}$  has multiplicity  $t > 1$ , necessitating nonsmooth minimization techniques. Under reasonable conditions, the set of points  $x \in \mathbf{R}^n$  for which  $\lambda_{\max}$  has fixed multiplicity  $t$  is a manifold of codimension  $\frac{t(t+1)}{2}$ , relative to which  $f$  is partly smooth [16].

For illustration, in Figure 5 we show convergence of the bundle method, BFGS, and bundle Newton method on a typical trial for this problem using random data. All algorithms were run without termination conditions until numerical issues prevented any further progress. In this example for  $n = 50$  matrices in  $\mathbf{R}^{25 \times 25}$ , the optimal eigenvalue multiplicity was 6, and we again observe fast convergence of the bundle Newton method once the subdifferential dimension  $\frac{t(t+1)}{2} - 1 = 20$  can be identified.

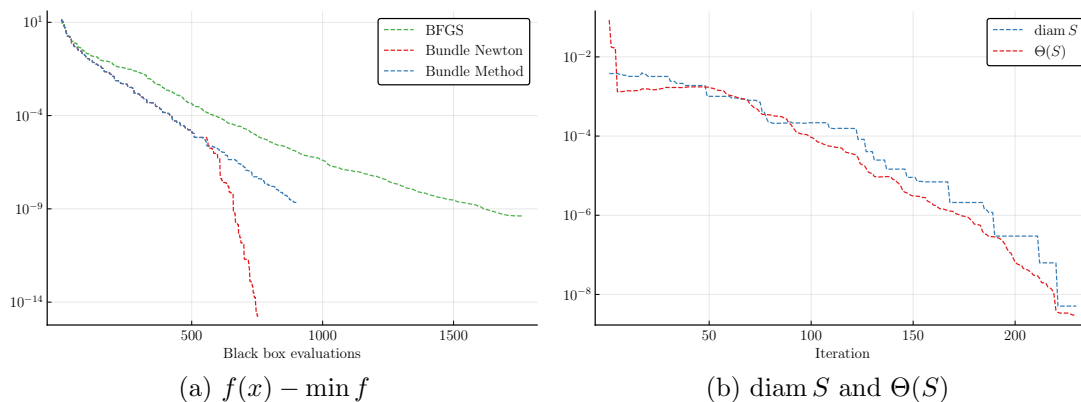


Figure 5: Function value convergence and optimality measures for the maximum eigenvalue function (8.7) for  $n = 50$  symmetric matrices in  $\mathbf{R}^{25 \times 25}$ .

(Note that since the optimal objective value is unknown, we instead used the best value found after running the algorithms with a large number of random starting points. This introduces a slight bias in the accuracy reported for the bundle Newton method.) In Figure 6, we observe that the bundle Newton method achieves an eigenvalue clustering several orders of magnitude better than is possible with a bundle method or BFGS.

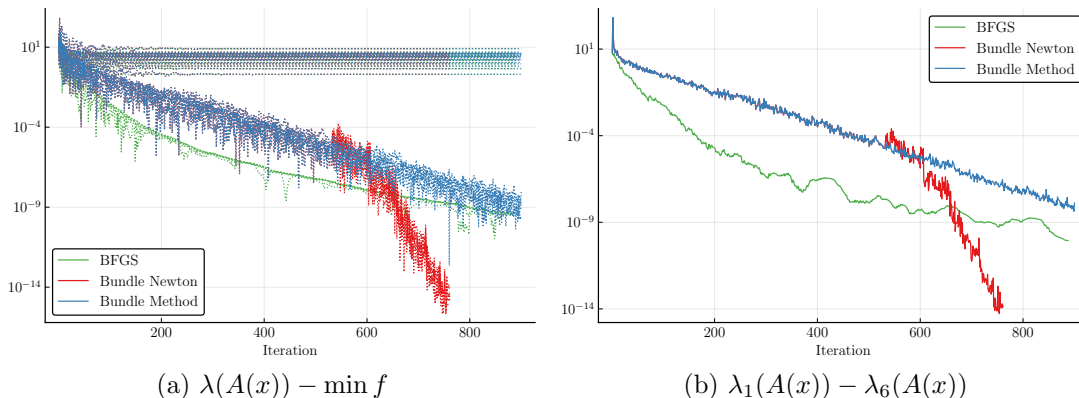


Figure 6: Clustering of the six largest eigenvalues of  $A(x)$ .

Using the active manifold to accelerate eigenvalue optimization is not new [29,30,35]. What is remarkable is that the bundle Newton method, combined with a first phase such as a traditional bundle method, rapidly converges to the minimizer *without* any structural knowledge of the function.

### First-order analogues

The Newton philosophy that we explore in this work is suggestive even in the more usual setting where Hessians are unavailable. One straightforward first-order analogue of Algorithm 2,2 replaces the Hessians by suitably tuned multiples of the identity matrix. Simple implementations seem effective on max functions: a broader investigation is the topic of ongoing work.

## References

- [1] D.S. Atkinson and P.M. Vaidya. A cutting plane algorithm for convex programming that uses analytic centers. *Math. Programming*, 69:1–43, 1995.
- [2] S. Bubeck. Convex optimization: algorithms and complexity. *Foundations and Trends in Machine Learning*, 8:231–357, 2015.
- [3] J.V. Burke, A.S. Lewis, and M.L. Overton. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM J. Optim.*, 15:751–779, 2005.
- [4] F.H. Clarke. *Optimization and Nonsmooth Analysis*. Wiley Interscience, New York, 1983.

- [5] M. Coste. *An Introduction to Semialgebraic Geometry*. RAAG Notes, 78 pages, Institut de Recherche Mathématiques de Rennes, October 2002.
- [6] Y. Du and A. Ruszczyński. Rate of convergence of the bundle method. *J. Optim. Theory Appl.*, 173:908–922, 2017.
- [7] J.-L. Goffin and J.-P. Vial. On the computation of weighted analytic centers and dual ellipsoids with the projective algorithm. *Math. Programming*, 60:81–92, 1993.
- [8] G.H. Golub and C.F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, MD, fourth edition, 2013.
- [9] K. C. Kiwiel. Efficiency of proximal bundle methods. *J. Optim. Theory Appl.*, 104:589–603, 2000.
- [10] K.C. Kiwiel. *Methods of Descent for Nondifferentiable Optimization*, volume 1133 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1985.
- [11] Yin Tat Lee, A. Sidford, and Sam Chiu-wai Wong. A faster cutting plane method and its implications for combinatorial and convex optimization. In *FOCS 2015*, pages 1049–1065. IEEE Computer Soc., Los Alamitos, CA, 2015.
- [12] C. Lemaréchal. An extension of Davidon methods to non differentiable problems. *Math. Programming Stud.*, 3:95–109, 1975.
- [13] C. Lemaréchal, A. Nemirovskii, and Y. Nesterov. New variants of bundle methods. *Math. Programming*, 69:111–147, 1995.
- [14] C. Lemaréchal, F. Oustry, and C. Sagastizábal. The U-lagrangian of a convex function. *Transactions of the American Mathematical Society*, 352:711–729, 2000.
- [15] C. Lemaréchal and C. Sagastizábal. Practical aspects of the Moreau-Yosida regularization: theoretical preliminaries. *SIAM J. Optim.*, 7:367–385, 1997.
- [16] A.S. Lewis. Active sets, nonsmoothness, and sensitivity. *SIAM J. Optim.*, 13:702–725, 2002.
- [17] A.S. Lewis and M.L. Overton. Nonsmooth optimization via quasi-Newton methods. *Math. Program.*, 141:135–163, 2013.
- [18] L. Lukšan and J. Vlček. A bundle-Newton method for nonsmooth unconstrained minimization. *Math. Programming*, 83:373–391, 1998.
- [19] R. Mifflin. An algorithm for constrained optimization with semismooth functions. *Math. Oper. Res.*, 2:191–207, 1977.



- [20] R. Mifflin and C. Sagastizábal. Proximal points are on the fast track. *Journal of Convex Analysis*, 9:563–579, 2002.
- [21] R. Mifflin and C. Sagastizábal. A  $\mathcal{VU}$ -algorithm for convex minimization. *Math. Program.*, 104:583–608, 2005.
- [22] S.A. Miller and J. Malick. Newton methods for nonsmooth convex minimization: connections among  $\mathcal{U}$ -Lagrangian, Riemannian Newton and SQP methods. *Math. Program.*, 104:609–633, 2005.
- [23] A. S. Nemirovsky and D. B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. John Wiley & Sons, Inc., New York, 1983. Translated from the Russian and with a preface by E. R. Dawson.
- [24] Y. Nesterov. Complexity estimates of some cutting plane methods based on the analytic barrier. *Math. Programming*, 69:149–176, 1995.
- [25] Y. Nesterov. *Introductory Lectures on Convex Optimization*. Kluwer Academic, Dordrecht, The Netherlands, 2004.
- [26] Y. Nesterov. Smooth minimization of non-smooth functions. *Math. Program.*, 103:127–152, 2005.
- [27] Y. Nesterov and A. Nemirovskii. *Interior-Point Polynomial Algorithms in Convex Programming*. SIAM, Philadelphia, PA, 1994.
- [28] J. Nocedal and S.J. Wright. *Numerical Optimization*. Springer, New York, second edition, 2006.
- [29] F. Oustry. The  $\mathcal{U}$ -Lagrangian of the maximum eigenvalue function. *SIAM J. Optim.*, 9:526–549, 1999.
- [30] M.L. Overton. On minimizing the maximum eigenvalue of a symmetric matrix. In *Linear algebra in signals, systems, and control*, pages 150–169. SIAM, Philadelphia, PA, 1988.
- [31] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1:123–231, 2013.
- [32] S.M. Robinson. Perturbed Kuhn-Tucker points and rates of convergence for a class of nonlinear-programming algorithms. *Math. Programming*, 7:1–16, 1974.
- [33] E.K. Ryu and S. Boyd. A primer on monotone operator methods. *Appl. Comput. Math.*, 15:3–43, 2016.

- [34] C. Sagastizábal. A  $\mathcal{VU}$ -point of view of nonsmooth optimization. In *Proceedings of the International Congress of Mathematicians, Rio de Janeiro*, volume 3, pages 3785–3806, 2018.
- [35] A. Shapiro and M.K.H. Fan. On eigenvalue optimization. *SIAM J. Optim.*, 5:552–569, 1995.
- [36] P. M. Vaidya. A new algorithm for minimizing convex functions over convex sets. *Math. Programming*, 73:291–341, 1996.