

LIMITED MEMORY DISCRETE GRADIENT BUNDLE METHOD FOR NONSMOOTH DERIVATIVE FREE OPTIMIZATION

N. KARMITSA¹ AND A.M. BAGIROV²

Abstract: Typically, practical nonsmooth optimization problems involve functions with hundreds of variables. Moreover, there are many practical problems where the computation of even one subgradient is either a difficult or an impossible task. In such cases derivative free methods are the better (or only) choice since they do not use explicit computation of subgradients. However, these methods require a large number of function evaluations even for moderately large problems. In this paper, we propose an efficient derivative free limited memory discrete gradient bundle method for nonsmooth, possibly nonconvex optimization. The convergence of the proposed method is proved for locally Lipschitz continuous functions and the numerical experiments to be presented confirm the usability of the method especially for medium size and large-scale problems.

Keywords: Nondifferentiable optimization, derivative free optimization, limited memory methods, bundle methods, discrete gradient.

1 Introduction

We introduce a new derivative free method for solving unconstrained minimization problems of the form

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathbb{R}^n, \end{cases} \quad (\text{P})$$

where the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is supposed to be locally Lipschitz continuous (LLC). In particular, our aim is to design a derivative free method for solving Problem (P) with the size from medium to large (i.e. problems with more than 50 variables). Note that no differentiability or convexity assumptions are made and we do not assume that a (sub)gradient of the objective function can be evaluated.

Nonsmooth optimization (NSO) problems (P) are encountered in many application areas: for instance, in economics [48], mechanics [46], engineering [45], control theory [16], optimal shape design [27], machine learning [30], and data mining [5, 12] including cluster analysis [18] and classification [2, 3, 10]. Most of these problems are large-scale.

NSO problems are in general difficult to solve. The direct application of gradient-based methods to nonsmooth problems may lead to a failure in convergence, in optimality conditions, or in gradient approximation (see, e.g. [37]).

1. Department of Mathematics and Statistics, University of Turku, FI-20014 Turku, Finland. E-mail: napsu@karmitsa.fi

2. Centre for Informatics and Applied Optimization, School of Science, Information Technology and Engineering, University of Ballarat, University Drive, Mount Helen, PO Box 663, Ballarat, VIC 3353, Australia. E-mail: a.bagirov@ballarat.edu.au

Methods for solving NSO problems include subgradient methods (see e.g. [8, 9, 52, 53]), bundle methods (see e.g. [20, 28, 35, 39, 41, 42, 50, 51]), algorithms based on smoothing techniques [49], and the gradient sampling methods [13]. In most of these algorithms the computation of at least one arbitrary subgradient (generalized gradient [15]) at each iteration is required. However, there are many practical problems where computation of even one subgradient is a difficult (or an impossible) task. In such cases derivative free methods are the preferable (or the only) choice since they do not require explicit computation of subgradients.

Most of existing derivative free methods like the genetic algorithm (see, e.g. [21]) or Powell’s method (see, e.g. [19]) are inefficient for solving nonsmooth optimization problems with even tens of variables. For the other derivative free methods the convergence can only be proved under differentiability or strict differentiability assumptions (see, e.g. [4, 17, 19]). The exceptions are the discrete gradient method [7] and the approximative subgradient method [6], where the convergence are proved for quasidifferentiable semismooth functions.

In addition to problematic subgradients, nonconvexity adds another challenge; NSO is traditionally based on convex analysis and most solution methods rely strongly on the convexity of the problem. Fortunately, several nonconvex algorithms have been introduced only recently, for instance, in [1, 13, 25, 26, 33, 47]. Nevertheless, these algorithms need to compute subgradients at each iteration point, just like the convex solution methods.

In this paper we introduce a new limited memory discrete gradient bundle method for derivative free nonconvex nonsmooth optimization. The idea of the new method is to combine the discrete gradient method DGM [7] and the limited memory bundle method LMBM [23, 24, 25]. DGM is a derivative free method for (small-scale) nonconvex nonsmooth optimization while LMBM utilizes subgradient information to solve large-scale nonsmooth problems. Although these two are totally different methods they have some similarities in their structures that makes combining of them rather an interesting task. The convergence of the proposed method is proved for semismooth functions.

The paper is organized as follows. In Section 2 we recall some basic definitions and results from nonsmooth analysis. We also discuss briefly DGM and LMBM. In Section 3 we introduce the new method and study its convergence properties. The results of the numerical experiments are presented and discussed in Section 4 and, finally, Section 5 concludes the paper.

2 Background

In this section, we first recall some basic definitions and results from nonsmooth analysis. Then, we discuss basic ideas of LMBM and DGM.

2.1 Preliminaries

We denote by $\|\cdot\|$ the Euclidean norm in \mathbb{R}^n and by $\mathbf{a}^T \mathbf{b}$ the inner product of vectors \mathbf{a} and \mathbf{b} (bolded symbols are used for vectors). An open (closed) ball with center $\mathbf{x} \in \mathbb{R}^n$ and radius $r > 0$ is denoted by $B(\mathbf{x}; r)$ ($\bar{B}(\mathbf{x}; r)$). That is,

$$B(\mathbf{x}; r) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y} - \mathbf{x}\| < r\} \quad \text{and} \quad \bar{B}(\mathbf{x}; r) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y} - \mathbf{x}\| \leq r\}.$$

The *subdifferential* $\partial f(\mathbf{x})$ [15] of a LLC function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at any point $\mathbf{x} \in \mathbb{R}^n$ is given by

$$\partial f(\mathbf{x}) = \text{conv}\left\{ \lim_{i \rightarrow \infty} \nabla f(\mathbf{x}_i) \mid \mathbf{x}_i \rightarrow \mathbf{x} \text{ and } \nabla f(\mathbf{x}_i) \text{ exists} \right\},$$

where “conv” denotes the convex hull of a set. A vector $\boldsymbol{\xi} \in \partial f(\mathbf{x})$ is called a *subgradient*. The *Goldstein ε -subdifferential* [22]) with some $\varepsilon \geq 0$ is a set

$$\partial_\varepsilon^G f(\mathbf{x}) = \text{conv}\{ \partial f(\mathbf{y}) \mid \mathbf{y} \in \bar{B}(\mathbf{x}; \varepsilon) \}.$$

Note that $\partial f(\mathbf{x}) \subseteq \partial_\varepsilon^G f(\mathbf{x})$ for all $\varepsilon \geq 0$.

The point $\mathbf{x}^* \in \mathbb{R}^n$ is called *substationary* if $\mathbf{0} \in \partial f(\mathbf{x}^*)$. Substationarity is a necessary condition for local optimality and, in the convex case, it is also sufficient for global optimality. An optimization method is said to be *globally convergent* if starting from any arbitrary point \mathbf{x}_1 it generates a sequence $\{\mathbf{x}_k\}$ that converges to a substationary point \mathbf{x}^* , that is, $\{\mathbf{x}_k\} \rightarrow \mathbf{x}^*$ whenever $k \rightarrow \infty$.

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called *semismooth* at $\mathbf{x} \in \mathbb{R}^n$ if it is LLC at \mathbf{x} and for every $\mathbf{d} \in \mathbb{R}^n$ the limit

$$\lim_{\substack{\boldsymbol{\xi} \in \partial f(\mathbf{x} + h\mathbf{d}'), \\ \mathbf{d}' \rightarrow \mathbf{d}, h \rightarrow 0^+}} \boldsymbol{\xi}^T \mathbf{d}$$

exists [43]. If the function f is semismooth, the classical directional derivative $f'(\mathbf{x}, \mathbf{d}) = \lim_{h \rightarrow 0^+} h^{-1}[f(\mathbf{x} + h\mathbf{d}) - f(\mathbf{x})]$ exists and

$$f'(\mathbf{x}, \mathbf{d}) = \lim_{\substack{\boldsymbol{\xi} \in \partial f(\mathbf{x} + h\mathbf{d}'), \\ \mathbf{d}' \rightarrow \mathbf{d}, h \rightarrow 0^+}} \boldsymbol{\xi}^T \mathbf{d}.$$

The class of semismooth functions contains convex, concave, max-type and min-type functions inter alia.

Let us denote by

$$S_1 = \{\mathbf{g} \in \mathbb{R}^n \mid \|\mathbf{g}\| = 1\}$$

the sphere of the unit ball and by

$$P = \{z \mid z : \mathbb{R}_+ \rightarrow \mathbb{R}_+, \zeta > 0, \zeta^{-1}z(\zeta) \rightarrow 0, \zeta \rightarrow 0\}$$

the set of univariate positive infinitesimal functions. In addition, let

$$G = \{\mathbf{e} \in \mathbb{R}^n \mid \mathbf{e} = (e_1, \dots, e_n), |e_j| = 1, j = 1, \dots, n\}$$

be a set of all vertices of the unit hypercube in \mathbb{R}^n . We take any $\mathbf{g} \in S_1$, $\mathbf{e} \in G$, $z \in P$, a positive number $\alpha \in (0, 1]$, and compute $i = \text{argmax}\{|g_j| \mid j = 1, \dots, n\}$. For $\mathbf{e} \in G$ we define the sequence of n vectors $\mathbf{e}^j(\alpha) = (\alpha e_1, \alpha^2 e_2, \dots, \alpha^j e_j, 0, \dots, 0)$ $j = 1, \dots, n$. For $\mathbf{x} \in \mathbb{R}^n$ and $\zeta > 0$, consider the points

$$\mathbf{x}_0 = \mathbf{x} + \zeta \mathbf{g}, \quad \mathbf{x}_j = \mathbf{x}_0 + z(\zeta) \mathbf{e}^j(\alpha), \quad j = 1, \dots, n. \quad (1)$$

DEFINITION 2.1. Let $\mathbf{g} \in S_1$, $\mathbf{e} \in G$, $z \in P$, $\alpha \in (0, 1]$, $\zeta > 0$ and take $i = \text{argmax}\{|g_j| \mid j = 1, \dots, n\}$. The *discrete gradient* of the function $f : \mathbb{R}^n \rightarrow$

\mathbb{R} at the point $\mathbf{x} \in \mathbb{R}^n$ in the direction \mathbf{g} is the vector $\Gamma^i(\mathbf{x}, \mathbf{g}, \mathbf{e}, z, \zeta, \alpha) = (\Gamma_1^i, \dots, \Gamma_n^i) \in \mathbb{R}^n$ with the following coordinates:

$$\begin{aligned}\Gamma_j^i &= [z(\zeta)\alpha^j e_j]^{-1} [f(\mathbf{x}_j) - f(\mathbf{x}_{j-1})], \quad j = 1, \dots, n, \quad j \neq i, \\ \Gamma_i^i &= (\zeta g_i)^{-1} \left[f(\mathbf{x} + \zeta \mathbf{g}) - f(\mathbf{x}) - \zeta \sum_{j=1, j \neq i}^n \Gamma_j^i g_j \right].\end{aligned}$$

It follows from Definition 2.1 that

$$f(\mathbf{x} + \zeta \mathbf{g}) - f(\mathbf{x}) = \zeta \mathbf{g}^T \Gamma^i(\mathbf{x}, \mathbf{g}, \mathbf{e}, z, \zeta, \alpha) \quad (2)$$

for all $\mathbf{g} \in S_1$, $\mathbf{e} \in G$, $z \in P$, $\alpha \in (0, 1]$ and $\zeta > 0$.

REMARK 2.1. The discrete gradient is defined with respect to a given direction $\mathbf{g} \in S_1$. To compute the discrete gradient we first define a sequence of points $\mathbf{x}_0, \dots, \mathbf{x}_n$ (see (1)) and compute the values of the function f at these points. That is, we compute $n + 2$ values of f including the point \mathbf{x} . The i th coordinate is defined so that it satisfies equality (2) which can be considered as some version of the mean value theorem.

It has been proved in [7] that for any $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{g} \in S_1$, $\mathbf{e} \in G$, $\zeta > 0$, $z \in P$, and $\alpha > 0$ we have $\|\Gamma^i\| \leq C(n)L$, where $C(n) = (n^2 + 2n^{3/2} - 2n^{1/2})^{1/2}$ and L is a Lipschitz constant of function f at \mathbf{x} . Furthermore, the closed convex set of discrete gradients

$$\begin{aligned}V_0(\mathbf{x}, \zeta) &= \text{cl conv}\{\mathbf{v} \in \mathbb{R}^n \mid \exists (\mathbf{g} \in S_1, \mathbf{e} \in G, z \in P, \alpha > 0) \\ &\quad \text{such that } \mathbf{v} = \Gamma^i(\mathbf{x}, \mathbf{g}, \mathbf{e}, z, \zeta, \alpha)\}\end{aligned}$$

is an approximation to the subdifferential $\partial f(\mathbf{x})$ for sufficiently small $\zeta > 0$ as stated in the following proposition [7]:

PROPOSITION 2.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a semismooth function at \mathbf{x} . For $\zeta > 0$ and $\mathbf{g} \in S_1$ define*

$$o(\zeta, \mathbf{g}) = f(\mathbf{x} + \zeta \mathbf{g}) - f(\mathbf{x}) - \zeta f'(\mathbf{x}, \mathbf{g}).$$

If $\zeta^{-1}o(\zeta, \mathbf{g}) \rightarrow 0$ uniformly with respect to \mathbf{g} as $\zeta \rightarrow 0^+$, then for any $\varepsilon > 0$ there exists $\zeta_0 > 0$ such that

$$V_0(\mathbf{x}, \zeta) \subset \partial f(\mathbf{x}) + B(\mathbf{0}; \varepsilon)$$

for all $\zeta \in (0, \zeta_0)$.

2.2 Discrete Gradient Method

Next we briefly describe the discrete gradient method (DGM) by Bagirov et. al. More details can be found in [7]. The idea of DGM is to hybridize derivative free methods with bundle methods. In contrast with bundle methods, which require the computation of a single subgradient of the objective function at each trial point, DGM approximates subgradients by discrete gradients using function values only. Similarly to bundle methods the previous values of discrete gradients

are gathered into a bundle and the null step is used if the current search direction is not good enough.

We start by describing the direction finding procedure. As mentioned before, the closed convex set of discrete gradients $V_0(\mathbf{x}, \zeta)$ is an approximation to the subdifferential $\partial f(\mathbf{x})$ for sufficiently small $\zeta > 0$. Thus, it can be used to compute the descent direction for the objective. However, the computation of the whole set $V_0(\mathbf{x}, \zeta)$ is not easy, and therefore, in DGM only a few discrete gradients from the set are used to calculate the descent direction.

The DGM consists of an outer and an inner iterations. Let us denote by k the index of the outer iteration and by s the index of inner iteration. In what follows we use only the iteration counter s whenever possible without confusion. At every iteration k we first (i.e. $s = 1$) compute the discrete gradient $\mathbf{v}_1 = \Gamma^i(\mathbf{x}_{k_1}, \mathbf{g}_1, \mathbf{e}, z, \zeta, \alpha)$ (see Definition 2.1) with respect to any initial direction $\mathbf{g}_1 \in S_1$ and we set the initial bundle of discrete gradients $\bar{V}_1(\mathbf{x}_{k_1}) = \{\mathbf{v}_1\}$. Then we compute the vector

$$\bar{\mathbf{v}}_s = \operatorname{argmin}_{\mathbf{v} \in \bar{V}_s(\mathbf{x}_{k_s})} \|\mathbf{v}\|^2,$$

that is the distance between the convex hull $\bar{V}_s(\mathbf{x}_{k_s})$ of all computed discrete gradients and the origin. If this distance is less than a given tolerance $\delta > 0$ we accept the point \mathbf{x}_{k_s} as an approximate substationary point and go to the next outer iteration. Otherwise, we compute another search direction

$$\mathbf{g}_{s+1} = -\frac{\bar{\mathbf{v}}_s}{\|\bar{\mathbf{v}}_s\|}$$

and we check whether this is a descent direction. That is, we have

$$f(\mathbf{x} + \zeta \mathbf{g}_{s+1}) - f(\mathbf{x}) \leq -c_1 \zeta \|\bar{\mathbf{v}}_s\|,$$

with the given numbers $c_1 \in (0, 1)$ and $\zeta > 0$. In this case, we set $\mathbf{d}_{k_s} = \mathbf{g}_{s+1}$ and stop the direction finding procedure. Otherwise, we take a null step. That is, we compute another discrete gradient $\mathbf{v}_{s+1} = \Gamma^i(\mathbf{x}_{k_s}, \mathbf{g}_{s+1}, \mathbf{e}, z, \zeta, \alpha)$ into the direction \mathbf{g}_{s+1} , update the bundle of discrete gradients by

$$\bar{V}_{s+1}(\mathbf{x}_{k_s}) = \operatorname{conv}\{\bar{V}_s(\mathbf{x}_{k_s}) \cup \{\mathbf{v}_{s+1}\}\}$$

and continue the direction finding procedure with $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s}$ and $s = s + 1$. Note that, at the null steps the approximation of the subdifferential $\partial f(\mathbf{x})$ is improved. It has been proved in [7] that the direction finding procedure is terminating such that either the descent direction \mathbf{d}_{k_s} is found or $\|\bar{\mathbf{v}}_s\| \leq \delta$.

When the descent direction \mathbf{d}_{k_s} has been found, we need to compute the next (inner) iteration point $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s} + t_{k_s} \mathbf{d}_{k_s}$, where the step size t_{k_s} is defined as

$$t_{k_s} = \operatorname{argmax} \{t \geq 0 \mid f(\mathbf{x}_{k_s} + t \mathbf{d}_{k_s}) - f(\mathbf{x}_{k_s}) \leq -c_2 t \|\bar{\mathbf{v}}_{k_s}\|\}$$

with given $c_2 \in (0, c_1]$.

The pseudo-code of DGM is the following:

```

PROGRAM DGM
INITIALIZE  $\mathbf{x}_1 \in \mathbb{R}^n$  and  $k = 1$ ;
OUTER ITERATION
  Set  $s = 1$  and  $\mathbf{x}_{k_s} = \mathbf{x}_k$ ;
  Compute the first discrete gradient  $\mathbf{v}_{k_1}$ ;
  Set  $\bar{V}(\mathbf{x}_{k_1}) = \{\mathbf{v}_{k_1}\}$ ;
  WHILE the termination condition is not met
    INNER ITERATION
      Compute the vector  $\bar{\mathbf{v}}_{k_s} = \operatorname{argmin}_{\mathbf{v} \in \bar{V}(\mathbf{x}_{k_s})} \|\mathbf{v}\|^2$ ;
      IF  $\|\bar{\mathbf{v}}_{k_s}\| \leq \delta_k$  with  $\delta_k > 0$  s.t.  $\delta_k \searrow 0$  when  $k \rightarrow \infty$  THEN
        Set  $\mathbf{x}_{k+1} = \mathbf{x}_{k_s}$  and  $k = k + 1$ ;
        Go to the next OUTER ITERATION;
      ELSE
        Compute the search direction  $\mathbf{d}_{k_s} = -\bar{\mathbf{v}}_{k_s} / \|\bar{\mathbf{v}}_{k_s}\|$ ;
        Find a step size  $t^k$ ;
        IF Descent condition holds THEN
          Construct the following iteration  $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s} + t^k \mathbf{d}_{k_s}$ ;
          Compute a new discrete gradient  $\mathbf{v}_{k_{s+1}}$  at  $\mathbf{x}_{k_{s+1}}$ ;
          Set  $\bar{V}(\mathbf{x}_{k_{s+1}}) = \{\mathbf{v}_{k_{s+1}}\}$ ;
        ELSE
          Compute a new discrete gradient  $\mathbf{v}_{k_{s+1}}$  at  $\mathbf{x}_{k_s}$ 
            in direction  $\mathbf{d}_{k_s}$ ;
          Update the set  $\bar{V}(\mathbf{x}_{k_{s+1}}) = \operatorname{conv}\{\bar{V}(\mathbf{x}_{k_s}) \cup \mathbf{v}_{k_{s+1}}\}$ ;
          Set  $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s}$ ;
        END IF
      Set  $s = s + 1$  and go to the next INNER ITERATION;
    END IF
  END INNER ITERATION
END WHILE
END OUTER ITERATION
RETURN final solution  $\mathbf{x}_k$ ;
END DGM

```

It has been proved that DGM is globally convergent for quasidifferentiable LLC functions under the assumption that the set of discrete gradients uniformly approximates the subdifferential [7].

2.3 Limited Memory Bundle Method (LMBM)

In this subsection, we describe the limited memory bundle algorithm (LMBM) by Karmitsa (née Haarala) et. al. [23, 24, 25, 32] for solving general, possibly non-convex, large-scale NSO problems. Here we assume that at every point \mathbf{x} we can evaluate the value of the objective function $f(\mathbf{x})$ and one arbitrary subgradient $\boldsymbol{\xi}$ from the subdifferential $\partial f(\mathbf{x})$.

LMBM is a hybrid of the variable metric bundle methods [39, 54] and the limited memory variable metric methods (see e.g. [14]), where the first ones have been developed for small- and medium-scale nonsmooth optimization and the latter ones for smooth large-scale optimization. LMBM exploits the ideas of the variable metric bundle methods, namely the utilization of null steps, simple aggregation of subgradients, and the subgradient locality measures, but the search

direction \mathbf{d}_k is calculated using the limited memory approach. That is,

$$\mathbf{d}_k = -D_k \tilde{\boldsymbol{\xi}}_k,$$

where $\tilde{\boldsymbol{\xi}}_k$ is (an aggregate) subgradient and D_k is the limited memory variable metric update that, in the smooth case, represents the approximation of the inverse of the Hessian matrix. Note that the matrix D_k is not formed explicitly but the search direction \mathbf{d}_k is calculated using the limited memory approach.

In order to determine a new step into the search direction \mathbf{d}_k , LMBM uses so-called *line search procedure*: a new iteration point \mathbf{x}_{k+1} and a new auxiliary point \mathbf{y}_{k+1} are produced such that

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k + t_L^k \mathbf{d}_k & \text{and} \\ \mathbf{y}_{k+1} &= \mathbf{x}_k + t_R^k \mathbf{d}_k, & \text{for } k \geq 1 \end{aligned} \quad (3)$$

with $\mathbf{y}_1 = \mathbf{x}_1$, where $t_R^k \in (0, t_{max}]$ and $t_L^k \in [0, t_R^k]$ are step sizes, and $t_{max} > 1$ is the upper bound for the step size. A necessary condition for a *serious step* to be taken is to have

$$t_R^k = t_L^k > 0 \quad \text{and} \quad f(\mathbf{y}_{k+1}) \leq f(\mathbf{x}_k) - \varepsilon_L^k t_R^k w_k, \quad (4)$$

where $\varepsilon_L^k \in (0, 1/2)$ is a line search parameter and $w_k > 0$ represents the desirable amount of descent of f at \mathbf{x}_k . If the condition (4) is satisfied, we have $\mathbf{x}_{k+1} = \mathbf{y}_{k+1}$. On the other hand, a *null step* is taken if

$$t_R^k > t_L^k = 0 \quad \text{and} \quad -\beta_{k+1} + \mathbf{d}_k^T \boldsymbol{\xi}_{k+1} \geq -\varepsilon_R^k w_k, \quad (5)$$

where $\varepsilon_R^k \in (\varepsilon_L^k, 1/2)$ is a line search parameter, $\boldsymbol{\xi}_{k+1} \in \partial f(\mathbf{y}_{k+1})$, and β_{k+1} is the subgradient locality measure [38, 44] similar to standard bundle methods. That is,

$$\beta_{k+1} = \max\{|f(\mathbf{x}_k) - f(\mathbf{y}_{k+1}) + (\mathbf{y}_{k+1} - \mathbf{x}_k)^T \boldsymbol{\xi}_{k+1}|, \gamma \|\mathbf{y}_{k+1} - \mathbf{x}_k\|^2\}. \quad (6)$$

Here $\gamma \geq 0$ is a distance measure parameter supplied by the user. Parameter γ can be set to zero when f is convex.

In the case of a null step, we set $\mathbf{x}_{k+1} = \mathbf{x}_k$ but information about the objective function is increased because we store the auxiliary point \mathbf{y}_{k+1} and the corresponding auxiliary subgradient $\boldsymbol{\xi}_{k+1} \in \partial f(\mathbf{y}_{k+1})$.

LMBM uses the original subgradient $\boldsymbol{\xi}_k$ after the serious step and the aggregate subgradient $\tilde{\boldsymbol{\xi}}_k$ after the null step for direction finding (i.e. we set $\tilde{\boldsymbol{\xi}}_k = \boldsymbol{\xi}_k$ if the previous step was a serious step). The aggregation procedure is carried out by determining multipliers λ_i^k satisfying $\lambda_i^k \geq 0$ for all $i \in \{1, 2, 3\}$, and $\sum_{i=1}^3 \lambda_i^k = 1$ that minimize the function

$$\begin{aligned} \varphi(\lambda_1, \lambda_2, \lambda_3) &= [\lambda_1 \boldsymbol{\xi}_m + \lambda_2 \boldsymbol{\xi}_{k+1} + \lambda_3 \tilde{\boldsymbol{\xi}}_k]^T D_k [\lambda_1 \boldsymbol{\xi}_m + \lambda_2 \boldsymbol{\xi}_{k+1} + \lambda_3 \tilde{\boldsymbol{\xi}}_k] \\ &+ 2(\lambda_2 \beta_{k+1} + \lambda_3 \tilde{\beta}_k). \end{aligned} \quad (7)$$

Here $\boldsymbol{\xi}_m \in \partial f(\mathbf{x}_k)$ is the current subgradient (m denotes the index of the iteration after the latest serious step, i.e. $\mathbf{x}_k = \mathbf{x}_m$), $\boldsymbol{\xi}_{k+1} \in \partial f(\mathbf{y}_{k+1})$ is the auxiliary subgradient, and $\tilde{\boldsymbol{\xi}}_k$ is the current aggregate subgradient from the previous iteration ($\tilde{\boldsymbol{\xi}}_1 = \boldsymbol{\xi}_1$). In addition, β_{k+1} is the current subgradient locality measure and

$\tilde{\beta}_k$ is the current aggregate subgradient locality measure ($\tilde{\beta}_1 = 0$). The resulting aggregate subgradient $\tilde{\xi}_{k+1}$ and aggregate subgradient locality measure $\tilde{\beta}_{k+1}$ are computed by the formulae

$$\tilde{\xi}_{k+1} = \lambda_1^k \xi_m + \lambda_2^k \xi_{k+1} + \lambda_3^k \tilde{\xi}_k \quad \text{and} \quad \tilde{\beta}_{k+1} = \lambda_2^k \beta_{k+1} + \lambda_3^k \tilde{\beta}_k. \quad (8)$$

Due to this simple aggregation procedure only one trial point \mathbf{y}_{k+1} and the corresponding subgradient $\xi_{k+1} \in \partial f(\mathbf{y}_{k+1})$ need to be stored.

In LMBM both the limited memory BFGS (L-BFGS) and the limited memory SR1 (L-SR1) update formulae [14] are used in calculations of the search direction and the aggregate values. The idea of limited memory matrix updating is that instead of storing large $n \times n$ matrices D_k , one stores a certain (usually small) number of vectors obtained at the previous iterations of the algorithm, and uses these vectors to implicitly define the variable metric matrices. In the case of a null step, we use the L-SR1 update, since this update formula allows us to preserve the boundedness and some other properties of generated matrices which guarantee the global convergence of the method. Otherwise, since these properties are not required after a serious step, the more efficient L-BFGS update is employed (see, for more details, [23, 24, 25]).

As a stopping parameter, we use the value

$$w_k = -\tilde{\xi}_k^T \mathbf{d}_k + 2\tilde{\beta}_k$$

and we stop if $w_k \leq \varepsilon$ for some user specified $\varepsilon > 0$. The parameter w_k is also used during the line search procedure to represent the desirable amount of descent.

The pseudo-code of LMBM is the following:

```

PROGRAM LMBM
  INITIALIZE  $\mathbf{x}_1 \in \mathbb{R}^n$ ,  $\xi_1 \in \partial f(\mathbf{x}_1)$ , and  $\varepsilon > 0$ ;
  Set  $k = 1$  and  $\mathbf{d}_1 = -\xi_1$ ;
  WHILE the termination condition  $w_k \leq \varepsilon$  is not met
    Find step sizes  $t_L^k$  and  $t_R^k$ ;
    Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_L^k \mathbf{d}_k$ ;
    Evaluate  $f(\mathbf{x}_{k+1})$  and  $\xi_{k+1} \in \partial f(\mathbf{x}_k + t_R^k \mathbf{d}_k)$ ;
    IF  $t_L^k > 0$  THEN
      SERIOUS STEP
        Compute the search direction  $\mathbf{d}_{k+1}$  using  $\xi_{k+1}$  and L-BFGS
          update;
      END SERIOUS STEP
    ELSE
      NULL STEP
        Compute the aggregate subgradient  $\tilde{\xi}_{k+1}$ ;
        Compute the search direction  $\mathbf{d}_{k+1}$  using  $\tilde{\xi}_{k+1}$  and L-SR1
          update;
      END NULL STEP
    END IF
    Set  $k = k + 1$ ;
  END WHILE
  RETURN final solution  $\mathbf{x}_k$ ;
END LMBM

```

Under the upper semismoothness assumption [11] LMBM can be proved to be globally convergent for LLC objective functions [23, 25].

3 Limited Memory Discrete Gradient Bundle Method

In this section we introduce the new derivative free limited memory discrete gradient bundle method (LDGBM).

As mentioned in the introduction LMBM and DGM have some similarities in their structures. For instance, both of these methods wipe out the old information whenever the serious step occurs. This property is different from standard bundle methods (see e.g. [42]) where the old information is collected near the current iteration point and stored to be used in the next iterations nonetheless of the step in question. In practice, storing all the old information may have several disadvantages: first, it needs storage space, although unbounded storage requirement may be addressed by so-called aggregation procedure introduced in [34]; second, it adds computational costs; and, what is the worst, it may store and use information that is no longer relevant due to the fact that it might have been collected far away from the current iteration point. The last point may be especially problematic in nonconvex cases.

In LMBM we bundle the subgradients that are computed in a small neighborhood of the current iteration point. This is similar to standard bundle methods although in LMBM we use this information only after null steps and, at the most, three subgradients are needed. On the other hand, in DGM discrete gradients calculated only at the current iteration point with respect to different directions are gathered into a bundle (see Definition 2.1 and Subsection 2.2). In our new method we combine these ideas and compute discrete gradients in a small neighborhood of the current iteration point with respect to different directions.

In DGM, similar to standard bundle methods, a quadratic subproblem needs to be solved to find the discrete gradient with the shortest norm and, as a consequence, to calculate the search direction. Now, instead of bundling an unlimited number of discrete gradients in null steps and computing the shortest norm, we compute the convex combination of at most three discrete gradients and the search direction is calculated using limited memory approach. Thus, a possibly time consuming quadratic direction finding problem needs not to be solved and also the difficulty with the unbounded amount of storage needed in DGM has been dealt with.

The obvious difference between LDGBM and LMBM is that we now use discrete gradients instead of subgradients of the objective function. In addition, we use both inner and outer iterations in order to avoid too tight approximations to the subgradients at the beginning of computation (thus, we have a derivative free method). The inner iteration of LDGBM is essentially same as LMBM. That is, the search direction is computed by the formula

$$\mathbf{d}_{k_s} = -D_{k_s} \tilde{\mathbf{v}}_{k_s},$$

where s and k are the indices of inner and outer iterations, $\tilde{\mathbf{v}}_{k_s}$ is an aggregate discrete gradient and D_{k_s} is a limited memory variable metric update. In addition, the line search procedure (cf. (3) – (5)) is used to determine a new iteration and auxiliary points $\mathbf{x}_{k_{s+1}}$ and $\mathbf{y}_{k_{s+1}}$, and the aggregation procedure (cf. (7) and (8)) is used to compute a new aggregate discrete gradient $\tilde{\mathbf{v}}_{k_{s+1}}$ and a new aggregate subgradient locality measure $\tilde{\beta}_{k_{s+1}}$.

The first discrete gradient $\mathbf{v}_{1_1} = \Gamma^i(\mathbf{x}, \mathbf{g}_{1_1}, \mathbf{e}, z, \zeta, \alpha)$, where $i = \operatorname{argmax} \{|g_j| \mid j = 1, \dots, n\}$, is computed to an arbitrary initial direction $\mathbf{g}_{1_1} \in S_1$. After that

we always use the previous normalized search direction $\mathbf{g}_{k_{s+1}} = \mathbf{d}_{k_s} / \|\mathbf{d}_{k_s}\|$ to compute the next discrete gradient $\mathbf{v}_{k_{s+1}}$. Parameters $z \in P$, $\zeta > 0$, and $\alpha > 0$ are selected similarly to DGM [7].

The inner iteration is terminated if we have

$$\frac{1}{2} \|\tilde{\mathbf{v}}_{k_s}\|^2 + \tilde{\beta}_{k_s} \leq \delta_k$$

for some outer iteration parameter $\delta_k > 0$.

We use here a new adaptive updating strategy for the selection of outer iteration parameter δ_k . At the beginning, the outer iteration parameter δ_1 is set to a large number. Each time the inner iteration is terminated we set

$$\delta_{k+1} = \min\{\sigma\delta_k, w_{k_s}\},$$

where $\sigma \in (0, 1)$ and $w_{k_s} = -\tilde{\mathbf{v}}_{k_s}^T \mathbf{d}_{k_s} + 2\tilde{\beta}_{k_s}$. Similarly to LMBM, the parameter w_{k_s} is used also during the line search procedure to represent the desirable amount of descent (cf. (4) and (5)).

Let us assume that the sequences $z_k \in P$, $\zeta_k > 0$, $z_k \rightarrow 0^+$, $\zeta_k \rightarrow 0^+$, $k \rightarrow \infty$, a sufficiently small number $\alpha > 0$ and the line search parameters $\varepsilon_L^{k_s} \in (0, 1/2)$ and $\varepsilon_R^{k_s} \in (\varepsilon_L^{k_s}, 1/2)$ are given. The pseudo-code of LDGBM is the following (note that obviously more details are given here than before):

```

PROGRAM LDGBM
INITIALIZE  $\mathbf{x}_1 \in \mathbb{R}^n$ ,  $\mathbf{g}_{1_1} \in S_1$ ,  $\varepsilon > 0$ ,  $\delta_1 > \varepsilon$ ,  $\zeta_1 > 0$ ,  $\sigma, \epsilon \in (0, 1)$  and  $k = 1$ ;
OUTER ITERATION
  Set  $s = 1$  and  $\mathbf{x}_{k_1} = \mathbf{x}_k$ ;
  WHILE the termination condition  $\delta_k \leq \varepsilon$  is not met
    INNER ITERATION
      SERIOUS STEP 1
        Compute the discrete gradient  $\mathbf{v}_{k_s} \in V_0(\mathbf{x}_{k_s}, \zeta_k)$  in
          direction  $\mathbf{g}_{k_s}$ ;
        Set  $m = s$ ,  $\tilde{\mathbf{v}}_{k_s} = \mathbf{v}_{k_s}$ , and  $\tilde{\beta}_{k_s} = 0$ ;
        Compute the search direction  $\mathbf{d}_{k_s}$  using  $\tilde{\mathbf{v}}_{k_s}$  and L-BFGS
          update;
      END SERIOUS STEP 1
      INNER ITERATION TERMINATION
        IF  $1/2\|\tilde{\mathbf{v}}_{k_s}\|^2 + \tilde{\beta}_{k_s} \leq \delta_k$  THEN;
          Set  $\mathbf{x}_{k+1} = \mathbf{x}_{k_s}$ ,  $\mathbf{g}_{k+1_1} = \mathbf{d}_{k_s}/\|\mathbf{d}_{k_s}\|$ ,  $\zeta_{k+1} = \epsilon\zeta_k$ ,
             $\delta_{k+1} = \min\{\sigma\delta_k, -\tilde{\mathbf{v}}_{k_s}^T \mathbf{d}_{k_s} + 2\tilde{\beta}_{k_s}\}$ , and  $k = k + 1$ ;
          Go to the next OUTER ITERATION;
        END IF
      END INNER ITERATION TERMINATION
      Find step sizes  $t_L^{k_s}$  and  $t_R^{k_s}$ , and the subgradient
        locality measure  $\beta_{k_{s+1}}$ ;
      IF  $t_L^{k_s} > 0$  THEN
        SERIOUS STEP 2
          Construct the iteration  $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s} + t_L^{k_s} \mathbf{d}_{k_s}$ ;
          Set  $\mathbf{g}_{k_{s+1}} = \mathbf{d}_{k_s}/\|\mathbf{d}_{k_s}\|$ ;
          Set  $s = s + 1$  and go to the next SERIOUS STEP 1;
        END SERIOUS STEP 2
      ELSE
        NULL STEP
          Construct the trial point  $\mathbf{y}_{k_{s+1}} = \mathbf{x}_{k_s} + t_R^{k_s} \mathbf{d}_{k_s}$ ;
          Compute the new discrete gradient  $\mathbf{v}_{k_{s+1}} \in V_0(\mathbf{y}_{k_{s+1}}, \zeta_k)$ 
            at the point  $\mathbf{y}_{k_{s+1}}$  in the direction  $\mathbf{g}_{k_{s+1}} = \mathbf{d}_{k_s}/\|\mathbf{d}_{k_s}\|$ ;
          Compute the aggregate values
             $\tilde{\mathbf{v}}_{k_{s+1}} = \lambda_1^{k_s} \mathbf{v}_{k_m} + \lambda_2^{k_s} \mathbf{v}_{k_{s+1}} + \lambda_3^{k_s} \tilde{\mathbf{v}}_{k_s}$  and
             $\tilde{\beta}_{k_{s+1}} = \lambda_2^{k_s} \beta_{k_{s+1}} + \lambda_3^{k_s} \tilde{\beta}_{k_s}$ ;
          Compute the new search direction  $\mathbf{d}_{k_{s+1}}$  using  $\tilde{\mathbf{v}}_{k_{s+1}}$ 
            and L-SR1 update;
          Set  $\mathbf{x}_{k_{s+1}} = \mathbf{x}_{k_s}$  and  $s = s + 1$ ;
          Go to the INNER ITERATION TERMINATION;
        END NULL STEP
      END IF
    END INNER ITERATION
  END WHILE
END OUTER ITERATION
RETURN final solution  $\mathbf{x}_k$ ;
END LDGBM

```

The discrete gradient is computed according to Definition 2.1 and Remark 2.1. Similarly to LMBM the search direction and the aggregate values are computed by using the L-BFGS update after serious steps and L-SR1 update otherwise (see

[14] for more details of limited memory matrix updating and [24, 25] for their usage in LMBM). The step sizes $t_R^{k_s} \in (0, t_{max}]$ and $t_L^{k_s} \in [0, t_R^{k_s}]$ with $t_{max} > 1$ are computed such that either condition (4) for serious steps or condition (5) for null steps is satisfied. In addition, the subgradient locality measure β_{k_s+1} as well as the multipliers $\lambda_i^{k_s}$ satisfying $\lambda_i^{k_s} \geq 0$ for all $i \in \{1, 2, 3\}$, and $\sum_{i=1}^3 \lambda_i^{k_s} = 1$ utilized in the aggregation procedure are computed similarly to LMBM (see equations (6) and (7)).

3.1 Convergence Analysis

We now prove the global convergence of the LDGBM-algorithm. In addition to assuming that the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is LLC, the level set $\{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}_1)\}$ is supposed to be bounded for every starting point $\mathbf{x}_1 \in \mathbb{R}^n$. Furthermore, we assume that each execution of the line search procedure is finite and the assumptions of Proposition 2.2 are satisfied. The line search procedure has been proved to be finite under the assumption of semismoothness (see [54]). On the other hand, in Proposition 2.2 we assume that function f is semismooth. Since the class of semismooth functions includes the class of upper semismooth functions, we here assume that the objective function f is semismooth.

DEFINITION 3.1. Let $\varepsilon > 0$ and $\zeta > 0$. We define the ε -set of discrete gradients by

$$V_\varepsilon(\mathbf{x}, \zeta) = \text{conv}\{V_0(\mathbf{y}, \zeta) \mid \mathbf{y} \in \bar{B}(\mathbf{x}; \varepsilon)\}.$$

LEMMA 3.2. For all inner iterations k_s , there exists $\varepsilon_{k_s} \geq 0$ such that

$$\tilde{\mathbf{v}}_{k_s} \in V_{\varepsilon_{k_s}}(\mathbf{x}_{k_s}, \zeta_k).$$

PROOF. After a serious step we have $\tilde{\mathbf{v}}_{k_s+1} = \mathbf{v}_{k_s+1} \in V_0(\mathbf{x}_{k_s+1}, \zeta_k)$ and, hence, the case is trivial. After a null step the aggregate discrete gradient $\tilde{\mathbf{v}}_{k_s+1}$ is computed as a convex combination of three discrete gradient: $\mathbf{v}_{k_s} \in V_0(\mathbf{x}_{k_s}, \zeta_k)$, $\mathbf{v}_{k_s+1} \in V_0(\mathbf{y}_{k_s+1}, \zeta_k)$ and $\tilde{\mathbf{v}}_{k_s}$. For null steps we use the scaled direction vector $\theta_k \mathbf{d}_k$ with $\theta_k = \min\{1, C/\|\mathbf{d}_k\|\}$ and predefined $C > 0$ in the line search (for more details, see [25]). Thus, by (3) we have $\|\mathbf{y}_{k_s+1} - \mathbf{x}_{k_s}\| \leq t_{max}C$, and we have $\mathbf{v}_{k_s+1} \in V_{\varepsilon_{k_s}}(\mathbf{x}_{k_s}, \zeta_k)$ with some $0 < \varepsilon_{k_s} \leq t_{max}C$. Since the first aggregate discrete gradient $\tilde{\mathbf{v}}_{k_s}$ after a serious step is equal to \mathbf{v}_{k_s} , they both belong to the set $V_{\varepsilon_{k_s}}(\mathbf{x}_{k_s}, \zeta_k)$ with any $\varepsilon_{k_s} \geq 0$.

As a convex combination of these three discrete gradients, the second aggregate discrete gradient $\tilde{\mathbf{v}}_{k_s+1}$ after the serious step belongs to the set $V_{\varepsilon_{k_s}}(\mathbf{x}_{k_s}, \zeta_k)$ with some $0 < \varepsilon_{k_s} \leq t_{max}C$. The rest of the proof follows by induction. \square

LEMMA 3.3. Suppose that the assumptions of Proposition 2.2 are satisfied at every $\mathbf{y} \in \bar{B}(\mathbf{x}; \varepsilon)$ with some $\varepsilon > 0$. Then, there exists $\zeta_0 > 0$ such that the ε -set of discrete gradients approximates the Goldstein subdifferential by

$$V_\varepsilon(\mathbf{x}, \zeta) \subset \partial_\varepsilon^G f(\mathbf{x}) + B(\mathbf{0}; \varepsilon)$$

with all $\zeta \in (0, \zeta_0)$.

PROOF. Let $\mathbf{v} \in V_\varepsilon(\mathbf{x}, \zeta)$. Then $\mathbf{v} = \sum \lambda_i \mathbf{v}_i^{\mathbf{y}}$, where $\mathbf{v}_i^{\mathbf{y}} \in V_0(\mathbf{y}^i, \zeta)$, $\mathbf{y}^i \in \bar{B}(\mathbf{x}; \varepsilon)$, $\sum \lambda_i = 1$ and $\lambda_i \geq 0$ for all i . Now, by Proposition 2.2 there exists ζ_0^i for any

$\varepsilon > 0$ such that $V_0(\mathbf{y}^i, \zeta) \subset \partial f(\mathbf{y}^i) + B(\mathbf{0}; \varepsilon)$ for all $\zeta \in (0, \zeta_0^i)$. Thus, we have $\mathbf{v}_i^y = \boldsymbol{\xi}_i^y + \mathbf{s}_i^\varepsilon$, where $\boldsymbol{\xi}_i^y \in \partial f(\mathbf{y}^i)$ and $\mathbf{s}_i^\varepsilon \in B(\mathbf{0}; \varepsilon)$. Hence,

$$\begin{aligned} \mathbf{v} &= \sum \lambda_i (\boldsymbol{\xi}_i^y + \mathbf{s}_i^\varepsilon) \\ &= \sum \lambda_i \boldsymbol{\xi}_i^y + \sum \lambda_i \mathbf{s}_i^\varepsilon \\ &\subset \text{conv}\{\partial f(\mathbf{y}^i) \mid \mathbf{y}^i \in \bar{B}(\mathbf{x}; \varepsilon)\} + B(\mathbf{0}; \varepsilon) \\ &= \partial_\varepsilon^G f(\mathbf{x}) + B(\mathbf{0}; \varepsilon) \end{aligned}$$

with all $\zeta \in (0, \min_i \{\zeta_0^i\})$. \square

DEFINITION 3.4. A point \mathbf{x} is called a (ζ, δ) -stationary point of the function f , if

$$\min_{\mathbf{v} \in V_0(\mathbf{x}, \zeta)} \|\mathbf{v}\| \leq \delta.$$

THEOREM 3.5. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a LLC semismooth function. Then the number of inner iteration loops in the LDGBM-algorithm is finite and, at the termination, the point x_{k+1} is a $(\zeta_k, \sqrt{2\delta_k})$ -stationary point of the function f .

PROOF. The inner iteration loop in the LDGBM-algorithm is essentially the same as LMBM with subgradients replaced by discrete gradients. It has been shown in [25] that LMBM terminates after finite number of iterations, if the stopping parameter (cf. δ_k here) is greater than zero. The usage of discrete gradients does not alter this result. Thus, the inner iteration loop in the LDGBM-algorithm is terminating after finite number of steps.

At the termination we have

$$\delta_k \geq \frac{1}{2} \|\tilde{\mathbf{v}}_{k_s}\|^2 + \tilde{\beta}_{k_s} \geq \frac{1}{2} \|\tilde{\mathbf{v}}_{k_s}\|^2.$$

By Lemma 3.2 we have $\tilde{\mathbf{v}}_{k_s} \in V_\mu(\mathbf{x}_{k_s}, \zeta_k)$ with some $\mu \geq 0$ and we set $\mathbf{x}_{k+1} = \mathbf{x}_{k_s}$ at the termination. Thus, we have

$$\min_{\mathbf{v} \in V_0(\mathbf{x}_{k+1}, \zeta_k)} \|\mathbf{v}\|^2 \leq \|\tilde{\mathbf{v}}_{k_s}\|^2 \leq 2\delta_k.$$

Then by Definition 3.4, x_{k+1} is a $(\zeta_k, \sqrt{2\delta_k})$ -stationary point of the function f . \square

THEOREM 3.6. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a semismooth LLC function and suppose that the level set $\{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}_1)\}$ is bounded. Then every accumulation point of the sequence $\{\mathbf{x}_k\}$ generated by the LDGBM-algorithm is substationary for f .

PROOF. Since f is LLC and the level set is bounded, we have

$$f^* = \inf\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n\} > -\infty.$$

By Theorem 3.5 the inner iterations generate $(\zeta_k, \sqrt{2\delta_k})$ -stationary points for all $k \geq 1$. Thus, it follows from Definition 3.4 that

$$\min\{\|\mathbf{v}\| \mid \mathbf{v} \in V_0(\mathbf{x}_{k+1}, \zeta_k)\} \leq \sqrt{2\delta_k} \quad (9)$$

for any $k \geq 1$. Since $\{f(\mathbf{x}_k)\}$ is a nonincreasing sequence, \mathbf{x}_k belongs to the level set with all $k \geq 1$. Moreover, the boundedness of the level set implies that the

sequence $\{\mathbf{x}_k\}$ has at least one accumulation point. Let $\bar{\mathbf{x}}$ be an accumulation point of $\{\mathbf{x}_k\}$ and let $\mathbf{x}_{k_i} \rightarrow \bar{\mathbf{x}}$ when $i \rightarrow \infty$. Then we have from (9) that

$$\min\{\|\mathbf{v}\| \mid \mathbf{v} \in V_0(\mathbf{x}_{k_i}, \zeta_{k_i-1})\} \leq \sqrt{2\delta_{k_i-1}}. \quad (10)$$

Now, by Lemma 3.3 for any $\mathbf{y} \in \bar{B}(\bar{\mathbf{x}}; \mu)$ with $\mu > 0$ there exists $\zeta_0 > 0$ such that

$$V_0(\mathbf{y}, \zeta) \subseteq V_\mu(\bar{\mathbf{x}}, \zeta) \subseteq \partial_\mu^G f(\bar{\mathbf{x}}) + B(\mathbf{0}; \mu) \quad (11)$$

for all $\zeta \in (0, \zeta_0)$. Since \mathbf{x}_{k_i} converges to $\bar{\mathbf{x}}$ there exists $i_0 \geq 1$ such that $\mathbf{x}_{k_i} \in B(\bar{\mathbf{x}}; \mu)$ for all $i \geq i_0$. Moreover, since $\delta_k \rightarrow 0$ and $\zeta_k \rightarrow 0$ as $k \rightarrow \infty$, there exists $k_0 > 1$ such that $\delta_k \leq \varepsilon$ with some $\varepsilon > 0$ and $\zeta_k < \zeta_0$ for all $k > k_0$. Then there exists $l_0 \geq i_0$ such that $k_l \geq k_0 + 1$ for all $l \geq l_0$. Therefore, from (10) and (11) we obtain

$$\min\{\|\mathbf{v}\| \mid \mathbf{v} \in \partial_\mu^G f(\bar{\mathbf{x}})\} \leq \mu + \sqrt{2\varepsilon}.$$

Now, since $\mu, \varepsilon > 0$ are arbitrary and the subdifferential is upper semicontinuous (see, e.g. [42]), we have $\mathbf{0} \in \partial f(\bar{\mathbf{x}})$. \square

4 Numerical Experiments

In this section we compare the implementations of the methods described in the previous sections and also SolvOpt, an implementation of Shor's r -algorithm. The test set used in our experiments consists of classical academic nonsmooth minimization problems from the literature and their extensions.

4.1 Solvers

The tested optimization codes with references to more detailed descriptions of the methods and their implementations are presented in Table 1.

Table 1: Tested optimization codes

Software	Author(s)	Method	Reference
SolvOpt	Kuntsevich & Kappel	Shor's r -algorithm	[29, 36, 52]
LMBM	Karmitsa	Limited memory bundle	[24, 25]
DGM	Bagirov	Discrete gradient	[7]
LDGBM	Karmitsa	L-discrete gradient bundle	

A brief description of each software and the references from where the code can be downloaded are in order.

SolvOpt (Solver for local nonlinear optimization problems) is an implementation of Shor's r -algorithm. In **SolvOpt** one can select to use either original subgradients or difference approximations of them (i.e. the user does not have to code difference approximations but to select one parameter to do this automatically). In our experiments we have used difference approximations.

The MatLab, C and Fortran source codes for **SolvOpt** are available for downloading from <http://www.kfunigraz.ac.at/imawww/kuntsevich/solvopt/>. In our experiments we used **SolvOpt** v.1.1 HP-UX FORTRAN-90 sources.

LMBM is an implementation of the limited memory bundle method specifically developed for large-scale nonsmooth problems (see Subection 2.3). In our experiments we used the difference approximations to calculate the approximations of the subgradients. The Fortran 77 source code and the mex-driver (for MatLab users) are available for downloading from <http://napsu.karmita.fi/lmbm/>.

DGM is a discrete gradient solver (see Subection 2.2) for derivative free optimization. To apply DGM, one only needs to be able to compute at every point \mathbf{x} the value of the objective function. The subgradient will then be approximated by discrete gradients. The Fortran 77 source code of DGM is available for downloading from <http://napsu.karmita.fi/dgm/>.

LDGBM is an implementation of the limited memory discrete gradient bundle method introduced in this paper. The Fortran 95 source code of LDGBM is available for downloading from <http://napsu.karmita.fi/ldgbm/>.

The experiments were performed on an Intel[®] Core[™] 2 CPU 1.80GHz. To compile the codes, we used `gfortran-4.3`, the GNU Fortran compiler.

4.2 Test problems and parameters

As already said the test set used in our experiments consists of classical academic nonsmooth minimization problems from the literature and their extensions. That is, we have used the small-scale problems 3.1 – 3.20 of [40] (omitting problems 3.13 and 3.14 in which DGM converged to unbounded minima) and larger problems 1 – 10 introduced in [24] that can be formulated with any number of variables. We have used here 50, 200 and 1000 variables.

Note that in the computation of both the difference approximations and the discrete gradients more than n function evaluations are needed per iteration. Therefore, already a problem with 200 variables can be considered as large one for the solvers not using (sub)gradient information and that is outstandingly true when dealing with nonsmooth problems.

We say that a solver finds the solution with respect to a tolerance $\varepsilon > 0$ if

$$\frac{f_{best} - f_{opt}}{1 + |f_{opt}|} \leq \varepsilon,$$

where f_{best} is a solution obtained with the solver and f_{opt} is the best known (or optimal) solution. We have accepted the results small-scale problems ($n \leq 50$) with respect to the tolerance $\varepsilon = 5 \cdot 10^{-4}$. With larger problems ($n \geq 200$), we have accepted the results with the tolerance $\varepsilon = 10^{-3}$. To obtain comparable results the stopping parameters of the codes were tuned by the procedure similar to [31].

For DGM the maximum size of the bundle was set to $\min\{n+3, 100\}$. With LMBM and LDGBM the natural choice for the bundle size is two (the larger bundle may be used but it is only utilized in step size selection). For all other parameters we have used the default settings of the codes.

4.3 Results

The results are summarized in Tables 2 and 3. We have compared the efficiency of the solvers both in terms of the computational time (*cpu*) and the number of function evaluations (*nf*, evaluations for short). In the smaller problems ($n \leq 20$, see Table 2) we only give numbers of evaluations since there was not a big difference in the computational time of the different solvers.

Table 2: Summary of the results in terms of evaluations for small-scale problems

P	n	SolvOpt	LMBM	DGM	LDGBM
3.1	2	270	989	1332	267
3.2	2	376	4 370	451	244
3.3	2	106	172	411	296
3.4	2	93	503	435	298
3.5	2	134	1 105	338	454
3.6	2	94	3 513	584	279
3.7	2	120	232	300	193
3.8	2	254	fail	188	487
3.9	2	122	300	371	391
3.10	2	118	633	672	234
3.11	4	166	6 654	1 662	710
3.12	5	258	7 649	2 945	1 059
3.15	6	1 039	3 059	5 544	4 286
3.16	10	896	fail	4 929	1 655
3.17	10	fail	fail	4 567	fail
3.18	12	1 036	24 698	9 308	4 862
3.19	20	6 277	7 243	5 541	4 243
3.20	20	fail	102376	5720	fail

In small-scale problems `SolvOpt` usually used less evaluations than any of the other solvers tested including our new solver while `DGM` was the most reliable solver tested (see Table 2). However, `LDGBM` usually used less evaluations than `DGM` and `LMBM` and the robustness of the solver was similar to `SolvOpt` and `LMBM`. Thus, we can say that our new solver is comparable with the existing solvers in small-scale settings.

Although `SolvOpt` solved small-scale problems both efficiently and robustly it was not in its element when solving the medium- and large-scale problems (see Table 3). The number of evaluations were still usually the smallest ones but the robustness was very weak.

Also our new solver `LDGBM` solved only five (out of ten) problems with 1000 variables with the desired accuracy. However, the failures obtained with `LDGBM` were mostly inaccurate result: with the relaxed tolerance $\varepsilon = 10^{-2}$ `LDGBM` failed only in one problem (problem 2) while with `SolvOpt` the number of failures is still five even if the relaxed tolerance was used. This indicates that with `LDGBM` a better accuracy might be achieved with more tight stopping criterion.

In smaller problems ($n \leq 200$) the robustness of `LDGBM` was as good as that of `DGM` which, on the other hand, was the most robust of the existing solvers. `LDGBM` also used clearly less evaluations and computational time than `DGM` and `LMBM`.

Table 3: Summary of the results for medium and large-scale problems

P	n	SolvOpt <i>nf/cpu</i>	LMBM <i>nf/cpu</i>	DGM <i>nf/cpu</i>	LDGBM <i>nf/cpu</i>
1	50	32 195/0.17	54 564/5.20	27 025/0.03	26 284/0.05
2	50	9 984/0.98	fail	fail	fail
3	50	fail	114 671/0.26	103 200/1.35	12 588/0.04
4	50	fail	97 701/1.09	89 616/1.86	13 858/0.17
5	50	2 237/0.03	84 240/0.81	36 083/0.27	6 548/0.07
6	50	854/0.01	106 807/0.74	9 201/0.04	3 969/0.02
7	50	3 728/0.13	127 804/3.05	41 548/1.09	17 760/0.47
8	50	fail	fail	92 527/0.95	32 915/0.08
9	50	fail	79 547/0.13	26 112/0.03	4 163/0.01
10	50	fail	169 377/0.40	52 670/0.74	16 177/0.04
1	200	555 779/9.23	848 449/6.54	312 544/0.59	598 321/2.35
2	200	48 770/79.86	fail	379 523/313.80	fail
3	200	fail	1 165 657/10.05	553 314/14.10	173 305/1.68
4	200	fail	503 113/22.24	568 910/32.67	104 501/4.38
5	200	12 889/0.78	46 957/1.79	261 919/7.02	22 860/0.91
6	200	fail	488 559/12.88	fail	31 825/0.73
7	200	fail	418 623/39.47	378 468/37.66	68 331/7.14
8	200	fail	325 371/1.63	862 819/17.09	80 260/0.53
9	200	19 063/0.32	313 324/1.99	163 043/0.59	12 540/0.09
10	200	fail	546 940/4.75	628 636/20.48	69 233/0.69
1	1000	fail	fail	10 592 047/92.43	fail
2	1000	44 213/1812.13	fail	fail	fail
3	1000	fail	4 910 632/208.60	3 579 840/111.49	1 484 246/61.77
4	1000	fail	1 314 594/282.92	3 085 985/452.59	762 840/165.82
5	1000	44 206/14.64	317 113/58.52	1 593 326/197.64	545 135/105.97
6	1000	fail	4 001 432/511.81	12 583 470/1107.24	fail
7	1000	fail	fail	4 455 499/1609.18	fail
8	1000	fail	2 302 536/55.69	2 993 617/54.78	182 975/4.67
9	1000	fail	820 998/25.62	993 300/16.12	70 181/2.17
10	1000	fail	fail	fail	fail

5 Conclusions

In this paper, we have described a new limited memory discrete gradient bundle method (LDGBM) for unconstrained nonsmooth optimization. LDGBM is a derivative free method and thus applicable to a broad class of optimization problems even if the (sub)gradient is not available. We have proved the global convergence of the method for locally Lipschitz continuous semismooth objective functions, which are not necessarily differentiable or convex.

The numerical experiments reported confirm that LDGBM is efficient for both convex and nonconvex nonsmooth optimization problems. With small-scale problems ($n \leq 20$) LDGBM was comparable with the existing solvers and with large numbers of variables it usually used clearly less function evaluations than the other solvers tested.

We can conclude that LDGBM is a good alternative to existing derivative free nonsmooth optimization algorithms and for medium- and large-scale problems it might well be the best choice.

Acknowledgements

The work was financially supported by the University of Turku (Finland) and the University of Ballarat (Australia).

References

- [1] APKARIAN, P., NOLL, D., AND PROT, O. A trust region spectral bundle method for non-convex eigenvalue optimization. *SIAM Journal on Optimization* 19, 1 (2008), 281–306.
- [2] ASTORINO, A., AND FUDULI, A. Nonsmooth optimization techniques for semi-supervised classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29, 12 (2007), 2135–2142.
- [3] ASTORINO, A., FUDULI, A., AND GORGONE, E. Nonsmoothness in classification problems. *Optimization Methods and Software* 23, 5 (2008), 675–688.
- [4] AUDET, C., AND DENNIS, J. E. J. Analysis of generalized pattern searches. *SIAM Journal on Optimization* 13 (2003), 889–903.
- [5] ÄYRÄMÖ, S. *Knowledge Mining Using Robust Clustering*. PhD thesis, University of Jyväskylä, Department of Mathematical Information Technology, 2006.
- [6] BAGIROV, A. M., AND GANJEHLOU, A. N. An approximate subgradient algorithm for unconstrained nonsmooth, nonconvex optimization. *Mathematical Methods of Operations Research* 67 (2008), 187–206.
- [7] BAGIROV, A. M., KARASOZEN, B., AND SEZER, M. Discrete gradient method: A derivative free method for nonsmooth optimization. *Journal of Optimization Theory and Applications* 137 (2008), 317–334.
- [8] BECK, A., AND TEBOULLE, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters* 31, 3 (2003), 167–175.

- [9] BEN-TAL, A., AND NEMIROVSKI, A. Non-Euclidean restricted memory level method for large-scale convex optimization. *Mathematical Programming* 102, 3 (2005), 407–456.
- [10] BERGERON, C., MOORE, G., ZARETZKI, J., BRENNEMAN, C., AND BENNETT, K. Fast bundle algorithm for multiple instance learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Preprint, 2011.
- [11] BIHAIN, A. Optimization of upper semidifferentiable functions. *Journal of Optimization Theory and Applications* 4 (1984), 545–568.
- [12] BRADLEY, P. S., FAYYAD, U. M., AND MANGASARIAN, O. L. Mathematical programming for data mining: Formulations and challenges. *INFORMS Journal on Computing* 11 (1999), 217–238.
- [13] BURKE, J. V., LEWIS, A. S., AND OVERTON, M. L. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM Journal on Optimization* 15 (2005), 751–779.
- [14] BYRD, R. H., NOCEDAL, J., AND SCHNABEL, R. B. Representations of quasi-Newton matrices and their use in limited memory methods. *Mathematical Programming* 63 (1994), 129–156.
- [15] CLARKE, F. H. *Optimization and Nonsmooth Analysis*. Wiley-Interscience, New York, 1983.
- [16] CLARKE, F. H., LEDYAEV, Y. S., STERN, R. J., AND WOLENSKI, P. R. *Nonsmooth Analysis and Control Theory*. Springer, New York, 1998.
- [17] CUSTODIO, A. L., DENNIS, J. E., AND VICENTE, L. N. Using simplex gradients of nonsmooth functions in direct search methods. *IMA J. Numer. Anal.* 28, 4 (2008), 770–784.
- [18] DEMYANOV, V. F., BAGIROV, A. M., AND RUBINOV, A. M. A method of truncated codifferential with application to some problems of cluster analysis. *Journal of Global Optimization* 23, 1 (2002), 63–80.
- [19] FLETCHER, R. *Practical Methods of Optimization*, 2nd ed. John Wiley & Sons, Chichester, 1987.
- [20] GAUDIOSO, M., AND MONACO, M. F. Variants to the cutting plane approach for convex nondifferentiable optimization. *Optimization* 25 (1992), 65–75.
- [21] GOLDBERG, D. E. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley Publishing Company, Inc., Reading, MA, 1998.
- [22] GOLDSTEIN, A. A. Optimization of Lipschitz continuous functions. *Mathematical Programming* 13, 1 (1977), 14–22.
- [23] HAARALA, M. *Large-Scale Nonsmooth Optimization: Variable Metric Bundle Method with Limited Memory*. PhD thesis, University of Jyväskylä, Department of Mathematical Information Technology, 2004.
- [24] HAARALA, M., MIETTINEN, K., AND MÄKELÄ, M. M. New limited memory bundle method for large-scale nonsmooth optimization. *Optimization Methods and Software* 19, 6 (2004), 673–692.
- [25] HAARALA, N., MIETTINEN, K., AND MÄKELÄ, M. M. Globally convergent limited memory bundle method for large-scale nonsmooth optimization. *Mathematical Programming* 109, 1 (2007), 181–205.

- [26] HARE, W. L., AND SAGASTIZÁBAL, C. Redistributed proximal bundle method for nonconvex optimization. *SIAM Journal on Optimization* 20, 5 (2010), 2442–2473.
- [27] HASLINGER, J., AND NEITTAANMÄKI, P. *Finite Element Approximation for Optimal Shape, Material and Topology Design*, 2nd ed. John Wiley & Sons, Chichester, 1996.
- [28] HIRIART-URRUTY, J.-B., AND LEMARÉCHAL, C. *Convex Analysis and Minimization Algorithms II*. Springer-Verlag, Berlin, 1993.
- [29] KAPPEL, F., AND KUNTSEVICH, A. An implementation of Shor’s r -algorithm. *Computational Optimization and Applications* 15 (2000), 193–205.
- [30] KÄRKKÄINEN, T., AND HEIKKOLA, E. Robust formulations for training multilayer perceptrons. *Neural Computation* 16 (2004), 837–862.
- [31] KARMITSA, N., BAGIROV, A., AND MÄKELÄ, M. M. Comparing different nonsmooth optimization methods and software. *Optimization Methods and Software*, 27, 1 (2012), 131–153.
- [32] KARMITSA, N., MÄKELÄ, M. M., AND ALI, M. M. Limited memory interior point bundle method for large inequality constrained nonsmooth minimization. *Applied Mathematics and Computation* 198, 1 (2008), 382–400.
- [33] KARMITSA, N., TANAKA FILHO, M., AND HERSKOVITS, J. Globally Convergent Cutting Plane Method for Nonconvex Nonsmooth Minimization. *Journal of Optimization Theory and Applications* 148, 3 (2011), 528–549.
- [34] KIWIEL, K. C. *An aggregate subgradient method for nonsmooth convex minimization*. *Mathematical Programming* 27, (1983), 320–341.
- [35] KIWIEL, K. C. *Methods of Descent for Nondifferentiable Optimization*. Lecture Notes in Mathematics 1133. Springer-Verlag, Berlin, 1985.
- [36] KUNTSEVICH, A., AND KAPPEL, F. SolvOpt — the solver for local nonlinear optimization problems. Karl-Franzens University of Graz: Graz, Austria, 1997.
- [37] LEMARÉCHAL, C. Nondifferentiable optimization. In *Optimization*, G. L. Nemhauser, A. H. G. Rinnooy Kan, and M. J. Todd, Eds. Elsevier North-Holland, Inc., New York, 1989, pp. 529–572.
- [38] LEMARÉCHAL, C., STRODIOT, J.-J., AND BIHAIN, A. On a bundle algorithm for nonsmooth optimization. In *Nonlinear Programming*, O. L. Mangasarian, R. R. Mayer, and S. M. Robinson, Eds. Academic Press, New York, 1981, pp. 245–281.
- [39] LUKŠAN, L., AND VLČEK, J. Globally convergent variable metric method for convex nonsmooth unconstrained minimization. *Journal of Optimization Theory and Applications* 102, 3 (1999), 593–613.
- [40] LUKŠAN, L., AND VLČEK, J. Test problems for nonsmooth unconstrained and linearly constrained optimization. Technical Report 798, Institute of Computer Science, Academy of Sciences of the Czech Republic, Prague, 2000.
- [41] MÄKELÄ, M. M. Survey of bundle methods for nonsmooth optimization. *Optimization Methods and Software* 17, 1 (2002), 1–29.
- [42] MÄKELÄ, M. M., AND NEITTAANMÄKI, P. *Nonsmooth Optimization: Analysis and Algorithms with Applications to Optimal Control*. World Scientific Publishing Co., Singapore, 1992.

- [43] MIFFLIN, R. Semismooth and semiconvex functions in constrained optimization. *SIAM Journal on Control and Optimization* 15, 6 (1977), 959–972.
- [44] MIFFLIN, R. A modification and an extension of Lemaréchal’s algorithm for nonsmooth minimization. *Mathematical Programming Study* 17 (1982), 77–90.
- [45] MISTAKIDIS, E. S., AND STAVROULAKIS, G. E. *Nonconvex Optimization in Mechanics. Smooth and Nonsmooth Algorithms, Heuristics and Engineering Applications by the F.E.M.* Kluwert Academic Publishers, Dordrecht, 1998.
- [46] MOREAU, J. J., PANAGIOTOPOULOS, P. D., AND STRANG, G., Eds. *Topics in Nonsmooth Mechanics.* Birkhäuser Verlag, Basel, 1988.
- [47] NOLL, D., PROT, O., AND RONDEPIERRE, A. A proximity control algorithm to minimize non-smooth and non-convex functions. *Pacific Journal of Optimization* 4, 3 (2008), 569–602.
- [48] OUTRATA, J., KOČVARA, M., AND ZOWE, J. *Nonsmooth Approach to Optimization Problems With Equilibrium Constraints. Theory, Applications and Numerical Results.* Kluwert Academic Publisher, Dordrecht, 1998.
- [49] POLAK, E., AND O. R. J. Algorithms for finite and semi-finite min-max-min problems using adaptive smoothing techniques. *Journal of Optimization Theory and Applications* 119 (2003), 421–457.
- [50] SAGASTIZÁBAL, C., AND SOLODOV, M. An infeasible bundle method for nonsmooth convex constrained optimization without a penalty function or a filter. *SIAM Journal on Optimization* 16, 1 (2005), 146–169.
- [51] SCHRAMM, H., AND ZOWE, J. A version of the bundle idea for minimizing a nonsmooth function: Conceptual idea, convergence analysis, numerical results. *SIAM Journal on Optimization* 2, 1 (1992), 121–152.
- [52] SHOR, N. Z. *Minimization Methods for Non-Differentiable Functions.* Springer-Verlag, Berlin, 1985.
- [53] URYASEV, S. P. Algorithms for nondifferentiable optimization problems. *Journal of Optimization Theory and Applications* 71 (1991), 359–388.
- [54] VLČEK, J., AND LUKŠAN, L. Globally convergent variable metric method for nonconvex nondifferentiable unconstrained minimization. *Journal of Optimization Theory and Applications* 111, 2 (2001), 407–430.