

Using positive spanning sets to achieve d-stationarity with the Boosted DC Algorithm

F. J. Aragón Artacho* R. Campoy* P. T. Vuong†

February 13, 2020

*This paper is dedicated to Professor Marco A. López Cerdá
on the occasion of his 70th birthday*

Abstract

The Difference of Convex functions Algorithm (DCA) is widely used for minimizing the difference of two convex functions. A recently proposed accelerated version, termed BDCA for Boosted DC Algorithm, incorporates a line search step to achieve a larger decrease of the objective value at each iteration. Thanks to this step, BDCA usually converges much faster than DCA in practice. The solutions found by DCA are guaranteed to be critical points of the problem, but these may not be local minima. Although BDCA tends to improve the objective value of the solutions it finds, these are frequently just critical points as well. In this paper we combine BDCA with a simple Derivative-Free Optimization (DFO) algorithm to force the d-stationarity (lack of descent direction) at the point obtained. The potential of this approach is illustrated through some computational experiments on a Minimum-Sum-of-Squares clustering problem. Our numerical results demonstrate that the new method provides better solutions while still remains faster than DCA in the majority of test cases.

Keywords: Difference of convex functions; boosted difference of convex functions algorithm; positive spanning sets; d-stationary points; derivative-free optimization.

1 Introduction

In this paper, we are interested in solving the following unconstrained DC (difference of convex functions) optimization problem:

$$\min_{x \in \mathbb{R}^m} \{\phi(x) := g(x) - h(x)\} \quad (\mathcal{P})$$

*Department of Mathematics, University of Alicante, Alicante, Spain.

Email: francisco.aragon@ua.es and ruben.campoy@ua.es

†Department of Mathematics, University of Vienna, Austria and School of Mathematical Sciences, University of Southampton, SO17 1BJ, Southampton, UK.

Email: t.v.phan@soton.ac.uk

where $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ and $h : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper, closed and convex functions, and g is smooth, with the conventions:

$$\begin{aligned} (+\infty) - (+\infty) &= +\infty, \\ (+\infty) - \lambda &= +\infty \quad \text{and} \quad \lambda - (+\infty) = -\infty, \quad \forall \lambda \in]-\infty, +\infty[. \end{aligned}$$

Problem (\mathcal{P}) can be tackled by the well-known *DC algorithm (DCA)* [14, 15]. DC programming has become an active research field for the last few decades [9] and DCA has been successfully applied to many real-world problems arising in different fields (see, e.g., [10]). Although DCA performs well in practice, its convergence can be fairly slow for some particular problems. In order to speed up the scheme, an accelerated version of the algorithm, called *Boosted DC algorithm (BDCA)*, has been recently proposed in [2, 3]. The BDCA performs a line search at the point generated by the classical DCA, which allows to achieve a larger decrease in the objective value at each iteration. In the numerical experiments reported in [2, 3] it was shown that BDCA was not only faster than DCA, but also often found solutions with lower objective value. However, although both algorithms are proved to converge to critical points of (\mathcal{P}) , there is no guarantee that these points are local minima. For this reason, a simple trick to achieve better solutions consists in running the algorithms from different starting points. Another approach has been recently used in [13], where the authors incorporated an inertial term into the algorithm making it converge to *better* critical points. In the recent work [12], the author proposed a DC scheme which is able to compute d-stationary points. Although this algorithm permits to address problems where the function g is nonsmooth, the second component function h needs to be the pointwise maximum of finitely many differentiable functions.

The aim of this paper is to show that it is possible to combine BDCA with a simple DFO (Derivative-Free Optimization) routine to guarantee d-stationarity at the limit point obtained by the algorithm. As a representative application, we perform a set of numerical experiments on the Minimum Sum-of-Squares Clustering problem studied in [3] to illustrate this observation. This problem has many critical points, where both DCA and BDCA tend to easily get trapped in. As a byproduct of the DFO step, we observe that in some problems a single run of the new algorithm is able to provide better solutions than those obtained by multiple restarts of DCA.

The rest of this paper is organized as follows. In Section 2 we recall some preliminary results. We propose a new variant of BDCA, named BDCA+, in Section 3. The results of some numerical experiments are presented in Section 4, where we compare the performance of DCA, BDCA and BDCA+ on several test cases. We finish with some conclusions in Section 5.

2 Preliminaries

Throughout this paper, $\langle x, y \rangle$ denotes the inner product of $x, y \in \mathbb{R}^m$, and $\|\cdot\|$ corresponds to the induced norm given by $\|x\| = \sqrt{\langle x, x \rangle}$. For any extended real-valued function $f : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$, the set $\text{dom } f := \{x \in \mathbb{R}^m \mid f(x) < +\infty\}$ denotes the (effective) *domain* of f . A function f is *proper* if its domain is nonempty. The function f is *coercive* if $f(x) \rightarrow +\infty$ whenever $\|x\| \rightarrow +\infty$, and it is said to be *convex* if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \quad \text{for all } x, y \in \mathbb{R}^m \text{ and } \lambda \in [0, 1].$$

Further, f is *strongly convex* with strong convexity parameter $\rho > 0$ if $f - \frac{\rho}{2}\|\cdot\|^2$ is convex, i.e., when

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\rho}{2}\lambda(1 - \lambda)\|x - y\|^2,$$

for all $x, y \in \mathbb{R}^m$ and $\lambda \in [0, 1]$. For any convex function f , the *subdifferential* of f at $x \in \mathbb{R}^m$ is the set

$$\partial f(x) := \{w \in \mathbb{R}^m \mid f(y) \geq f(x) + \langle w, y - x \rangle, \forall y \in \mathbb{R}^m\}.$$

If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$, where $\nabla f(x)$ denotes the *gradient* of f at x . The one-sided *directional derivative* of f at x with respect to the direction $d \in \mathbb{R}^m$ is defined by

$$f'(x; d) := \lim_{t \searrow 0} \frac{f(x + td) - f(x)}{t}.$$

Before going to the main contribution of this paper in Section 3, we state our assumptions imposed on (\mathcal{P}) . We also recall some preliminary notions and basic results which will be used in the sequel.

2.1 Basic Assumptions

Assumption 1. Both functions g and h in (\mathcal{P}) are strongly convex on their domain for the same strong convexity parameter $\rho > 0$.

Assumption 2. The function h is subdifferentiable at every point in $\text{dom } h$; that is, $\partial h(x) \neq \emptyset$ for all $x \in \text{dom } h$.

Assumption 3. The function g is continuously differentiable on an open set containing $\text{dom } h$ and

$$\phi^* := \inf_{x \in \mathbb{R}^m} \phi(x) > -\infty.$$

Assumption 1 is not restrictive, as one can always rewrite the objective function as $\phi = (g + q) - (h + q)$ for any strongly convex function q (e.g., $q = \frac{\rho}{2}\|\cdot\|^2$). Observe that Assumption 2 holds for all $x \in \text{ri dom } h$ (by [17, Theorem 23.4]). A key property for our method is the smoothness of g in Assumption 3, which cannot be in general omitted (see [3, Example 3.2]).

2.2 Optimality Conditions

Under Assumptions 2 and 3 the following well-known necessary condition for local optimality holds.

Fact 2.1 (First-order necessary optimality condition). *If $x^* \in \text{dom } \phi$ is a local minimizer of problem (\mathcal{P}) , then*

$$\partial h(x^*) = \{\nabla g(x^*)\}. \tag{1}$$

Proof. See [16, Theorem 3']. □

Any point satisfying condition (1) is called a *d(irectional)-stationary point* of (\mathcal{P}) . We say that x^* is a *critical point* of (\mathcal{P}) if

$$\nabla g(x^*) \in \partial h(x^*).$$

Clearly, d-stationary points are critical points, but the converse is not true in general (see, e.g., [4, Example 1]). In our setting, the notion of critical point coincides with that of *Clarke stationarity*, which requires that zero belongs to the Clarke subdifferential at x^* (see, e.g., [5, Proposition 2]). The next result establishes that the d-stationary points of (\mathcal{P}) are precisely those points for which the directional derivative is zero for every direction.

Proposition 2.1. A point $x^* \in \text{dom } \phi$ is a d-stationary point of (\mathcal{P}) if and only if

$$\phi'(x^*; d) = 0, \quad \text{for all } d \in \mathbb{R}^m. \quad (2)$$

Proof. If x^* is a d-stationary point of (\mathcal{P}) , then by [17, Theorem 25.1] we know that h is differentiable at x^* . Therefore, for any $d \in \mathbb{R}^m$, we have

$$\phi'(x^*; d) = \langle \nabla g(x^*), d \rangle - \langle \nabla h(x^*), d \rangle = 0.$$

For the converse implication, pick any $v \in \partial h(x^*) \neq \emptyset$ (by Assumption 2) and observe that, for any $d \in \mathbb{R}^m$, we have that

$$\begin{aligned} \phi'(x^*; d) &= g'(x^*; d) - h'(x^*; d) \\ &= \langle \nabla g(x^*), d \rangle - \lim_{t \searrow 0} \frac{h(x^* + td) - h(x^*)}{t} \\ &\leq \langle \nabla g(x^*) - v, d \rangle. \end{aligned}$$

Hence, if x^* satisfies (2), one must have

$$\langle \nabla g(x^*) - v, d \rangle \geq 0, \quad \text{for all } d \in \mathbb{R}^m,$$

which is equivalent to $\nabla g(x^*) - v = 0$. As v was arbitrarily chosen in $\partial h(x^*)$, we conclude that $\partial h(x^*) = \{\nabla g(x^*)\}$. \square

2.3 DCA and Boosted DCA

In this section, we recall the iterative procedure DCA and its accelerated extension, BDCA, for solving problem (\mathcal{P}) . The DCA iterates by solving a sequence of approximating convex subproblems, as described next in Algorithm 1.

The key feature that makes the DCA work, stated next in Fact 2.2(a), is that the solution of (\mathcal{P}_k) provides a decrease in the objective value of (\mathcal{P}) along the iterations. Actually, an analogous result holds for the dual problem, see [14, Theorem 3]. In [2], the authors showed that the direction generated by the iterates of DCA, namely $d_k := y_k - x_k$, provides a descent direction of the objective function at y_k when the functions g and h in (\mathcal{P}) are assumed to be smooth. This result was later generalized in [3] to the case where h satisfies Assumption 2. The following result collects these properties.

Fact 2.2. Let x_k and y_k be the sequences generated by Algorithm 1 and set $d_k := y_k - x_k$, for all $k \in \mathbb{N}$. Then the following statements hold:

- (a) $\phi(y_k) \leq \phi(x_k) - \rho \|d_k\|^2$;
- (b) $\phi'(y_k; d_k) \leq -\rho \|d_k\|^2$;

Algorithm 1: DCA (DC Algorithm)

Input: An initial point $x_0 \in \mathbb{R}^m$ and a desired tolerance $\varepsilon \geq 0$;

```
1 begin
2    $k \leftarrow 0$ ;
3   Select  $u_k \in \partial h(x_k)$  and compute the unique solution  $y_k$  of
      
$$\min_{x \in \mathbb{R}^m} \{\phi_k(x) := g(x) - \langle u_k, x \rangle\}; \quad (\mathcal{P}_k)$$

4   if  $\|y_k - x_k\| \leq \varepsilon$  then
5     | stop and return  $y_k$ ;
6   else
7     |  $x_{k+1} = y_k$ ;
8   end
9    $k \leftarrow k + 1$  and go to line 3;
10 end
```

(c) there exists some $\delta_k > 0$ such that

$$\phi(y_k + \lambda d_k) \leq \phi(y_k) - \alpha \lambda^2 \|d_k\|^2, \quad \text{for all } \lambda \in [0, \delta_k[.$$

Proof. See [3, Proposition 3.1]. □

Thanks to Fact 2.2, once y_k has been found by DCA, one can achieve a larger decrease in the objective value of (\mathcal{P}) by moving along the descent direction d_k . Indeed, observe that

$$\phi(y_k + \lambda d_k) \leq \phi(y_k) - \alpha \lambda^2 \|d_k\|^2 \leq \phi(x_k) - (\rho + \alpha \lambda^2) \|d_k\|^2, \quad \text{for all } \lambda \in [0, \delta_k[.$$

This fact is the main idea of the BDCA [2, 3], whose iteration is described next in Algorithm 2.

Algorithmically, the BDCA is nothing more than the classical DCA with a line search procedure using an Armijo type rule. Note that the backtracking step in Algorithm 2 (lines 6–9) terminates finitely thanks to Fact 2.2(c).

We state next the basic convergence results for the sequences generated by BDCA (for more, see [2, 3]). Observe that DCA can be seen as a particular case of BDCA if one sets $\bar{\lambda}_k = 0$, so the following result applies to both Algorithms 1 and 2.

Fact 2.3. *For any $x_0 \in \mathbb{R}^m$, either Algorithm 2 (BDCA) with $\varepsilon = 0$ returns a critical point of (\mathcal{P}) , or it generates an infinite sequence such that the following properties hold.*

- (a) $\{\phi(x_k)\}$ is monotonically decreasing and convergent to some ϕ^* .
- (b) Any limit point of $\{x_k\}$ is a critical point of (\mathcal{P}) . In addition, if ϕ is coercive then there exists a subsequence of $\{x_k\}$ which converges to a critical point of (\mathcal{P}) .
- (c) It holds that $\sum_{k=0}^{+\infty} \|d_k\|^2 < +\infty$. Furthermore, if there is some $\bar{\lambda}$ such that $\lambda_k \leq \bar{\lambda}$ for all $k \geq 0$, then $\sum_{k=0}^{+\infty} \|x_{k+1} - x_k\|^2 < +\infty$.

Proof. See [3, Theorem 3.6]. □

Algorithm 2: BDCA (Boosted DC Algorithm)

Input: An initial point $x_0 \in \mathbb{R}^m$ and a desired tolerance $\varepsilon \geq 0$. Choose some parameters $\alpha > 0$ and $\beta \in]0, 1[$;

```
1 begin
2    $k \leftarrow 0$ ;
3   Select  $u_k \in \partial h(x_k)$  and compute the unique solution  $y_k$  of
      
$$\min_{x \in \mathbb{R}^m} \{\phi_k(x) := g(x) - \langle u_k, x \rangle\}; \quad (\mathcal{P}_k)$$

4    $d_k \leftarrow y_k - x_k$ ;
5   if  $\|d_k\| > \varepsilon$  then
6     Choose any  $\bar{\lambda}_k \geq 0$  and set  $\lambda_k \leftarrow \bar{\lambda}_k$ ;
7     while  $\phi(y_k + \lambda_k d_k) > \phi(y_k) - \alpha \lambda_k^2 \|d_k\|^2$  do
8        $\lambda_k \leftarrow \beta \lambda_k$ ;
9     end
10     $x_{k+1} \leftarrow y_k + \lambda_k d_k$ ;
11  else
12    stop and return  $y_k$ ;
13  end
14   $k \leftarrow k + 1$  and go to line 3;
15 end
```

2.4 Positive Spanning Sets

Most directional direct search methods are based on the use of positive spanning sets (see, e.g., [1, Section 5.6.3] and [7, Chapter 7]). Let us recall this concept here.

Definition 2.1. We call *positive span* of a set of vectors $\{v_1, v_2, \dots, v_r\} \subset \mathbb{R}^m$ the convex cone generated by this set, i.e.,

$$\{v \in \mathbb{R}^m : v = \alpha_1 v_1 + \dots + \alpha_r v_r, \quad \alpha_i \geq 0, i = 1, 2, \dots, r\}.$$

A set of vectors in \mathbb{R}^m is said to be a *positive spanning set* if its positive span is the whole space \mathbb{R}^m . A set $\{v_1, v_2, \dots, v_r\}$ is said to be *positively dependent* if one of the vectors is in the positive span generated by the remaining vectors; otherwise, the set is called *positively independent*. A *positive basis* in \mathbb{R}^m is a positively independent set whose positive span is \mathbb{R}^m .

Three well-known examples of positive spanning sets are given next.

Example 2.1 (Positive basis). Let e_1, e_2, \dots, e_m be the unit vectors of the standard basis in \mathbb{R}^m . Then the following sets are positive basis in \mathbb{R}^m :

$$D_1 := \{\pm e_1, \pm e_2, \dots, \pm e_m\}, \quad (3a)$$

$$D_2 := \{e_1, e_2, \dots, e_m, -\sum_{i=1}^m e_i\}, \quad (3b)$$

$$D_3 := \left\{ v_1, v_2, \dots, v_m, v_{m+1} \in \mathbb{R}^m, \quad \text{with} \quad \begin{array}{l} v_i^T v_j = \frac{-1}{m}, \text{ if } i \neq j, \\ \|v_i\| = 1, i = 1, 2, \dots, m+1. \end{array} \right\}. \quad (3c)$$

A possible construction for D_3 is given in [7, Corollary 2.6].

Recall that the BDCA provides critical points of (\mathcal{P}) which are not necessarily d-stationary points (Fact 2.3). Theoretically, see [14, Section 3.3], if x^* is a critical point which is not d-stationary, one could restart BDCA by taking $x_0 := x^*$ and choose $y_0 \in \partial h(x_0) \setminus \{\nabla g(x_0)\}$. Nonetheless, observe that this is only applicable when the algorithm converges in a finite number of iterations to x^* , which does not happen very often in practice (except for polyhedral DC problems, where even a global solution can be effectively computed if h is a piecewise linear function with a reasonable small number of pieces, see [14, §4.2]). Because of that, our goal is to design a variant of BDCA that generates a sequence converging to a d-stationary point. The following key result, proved in [6, Theorem 3.1], asserts that using positive spanning sets one can escape from points which are not d-stationary. We include its short proof.

Fact 2.4. *Let $\{v_1, v_2, \dots, v_r\}$ be a positive spanning set of \mathbb{R}^m . A point $x^* \in \text{dom } \phi$ is a d-stationary point of (\mathcal{P}) if and only if*

$$\phi'(x^*; v_i) \geq 0, \quad \text{for all } i = 1, 2, \dots, r. \quad (4)$$

Proof. The direct implication is an immediate consequence of Proposition 2.1. For the reverse implication, pick any $x^* \in \text{dom } \phi$ verifying (4) and choose any $d \in \mathbb{R}^m$. Since $\{v_1, v_2, \dots, v_r\}$ is a positive spanning set, there are $\alpha_1, \alpha_2, \dots, \alpha_r \geq 0$ such that

$$d = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_r v_r.$$

According to [17, Theorem 23.1], we have that

$$h'(x^*; d) \leq \alpha_1 h'(x^*; v_1) + \dots + \alpha_r h'(x^*; v_r).$$

Hence, we obtain

$$\begin{aligned} \phi'(x^*; d) &= g'(x^*; d) - h'(x^*; d) \\ &= \langle \nabla g(x^*), \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_r v_r \rangle - h'(x^*; d) \\ &\geq \sum_{i=1}^r \alpha_i \langle \nabla g(x^*), v_i \rangle - \sum_{i=1}^r \alpha_i h'(x^*; v_i) \\ &= \sum_{i=1}^r \alpha_i \phi'(x^*; v_i) \geq 0. \end{aligned}$$

Since d was arbitrarily chosen, then (2) holds and x^* is a d-stationary point of (\mathcal{P}) . \square

3 Forcing BDCA to converge to d-stationary points

In this section we propose a new variant of BDCA to solve problem (\mathcal{P}) , called BDCA+. The idea is to combine BDCA with a basic DFO routine which uses positive spanning sets. The first scheme aims at achieving a fast minimization of the objective function ϕ , while the second one is used to avoid converging to critical points for which there is at least a descent direction (i.e., they are not d-stationary points and, thus, they cannot be local minima). Let us make some comments about the new scheme BDCA+, which is stated in Algorithm 3.

- Subproblem (\mathcal{P}_k) in line 3 corresponds to the classical DCA step for solving (\mathcal{P}) .
- Lines 5 to 10 encode the boosting line search step used in BDCA. If the current iterate is (numerically) not a critical point, then the algorithm performs a line search step at y_k along the direction d_k to improve the objective values of (\mathcal{P}) .

- Line 11 to 19 correspond to a direct search DFO technique. It is run only when BDCA was stopped, in order to check if the point obtained is d-stationary. To this aim, it performs a backtracking search along each of the directions belonging to a positive spanning set D of \mathbb{R}^m . If it reaches a point whose objective value is smaller, then we move to that point and run BDCA again from there. Otherwise, there is not descent direction in D and, according to Fact 2.4, the point we have found must be (numerically) d-stationary.
- The choice $\bar{\lambda}_k = 0$ for all k is allowed, which corresponds to adding a direct search step to DCA.

Algorithm 3: BDCA+ (Boosted DC Algorithm combined with DFO)

Input: An initial point $x_0 \in \mathbb{R}^m$, a positive spanning set D . Choose three nonnegative parameters $\varepsilon_1, \eta, \tau \geq 0$, three positive ones $\alpha, \varepsilon_2, \bar{\mu} > 0$, and $\beta_1, \beta_2 \in]0, 1[$;

```

1 begin
2    $k \leftarrow 0, \mu \leftarrow \bar{\mu}$ ;
3   Select  $u_k \in \partial h(x_k)$  and compute the unique solution  $y_k$  of
                                     
$$\min_{x \in \mathbb{R}^m} \{\phi_k(x) := g(x) - \langle u_k, x \rangle\}; \quad (\mathcal{P}_k)$$

4    $d_k \leftarrow y_k - x_k$ ;
5   if  $\|d_k\| > \varepsilon_1$  then
6     Choose any  $\bar{\lambda}_k \geq 0$  and set  $\lambda_k \leftarrow \bar{\lambda}_k$ ;
7     while  $\phi(y_k + \lambda_k d_k) > \phi(y_k) - \alpha \lambda_k^2 \|d_k\|^2$  do
8        $\lambda_k \leftarrow \beta_1 \lambda_k$ ;
9     end
10     $x_{k+1} \leftarrow y_k + \lambda_k d_k$ ;
11  else
12     $\mu \leftarrow \eta \mu + \tau$ ;
13    if  $\phi(y_k + \mu v) < \phi(y_k)$  for some  $v \in D$  then
14       $x_{k+1} \leftarrow y_k + \mu v$ ;
15    else if  $\mu > \varepsilon_2$  then
16       $\mu \leftarrow \beta_2 \mu$  and go to line 12;
17    else
18      stop and return  $y_k$ ;
19    end
20   $k \leftarrow k + 1$  and go to line 3;
21 end

```

The following constructive example serves to illustrate the different behavior of DCA, BDCA and BDCA+.

Example 3.1 ([3, Example 3.3]). Consider the function $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$\phi(x, y) := x^2 + y^2 + x + y - |x| - |y|.$$

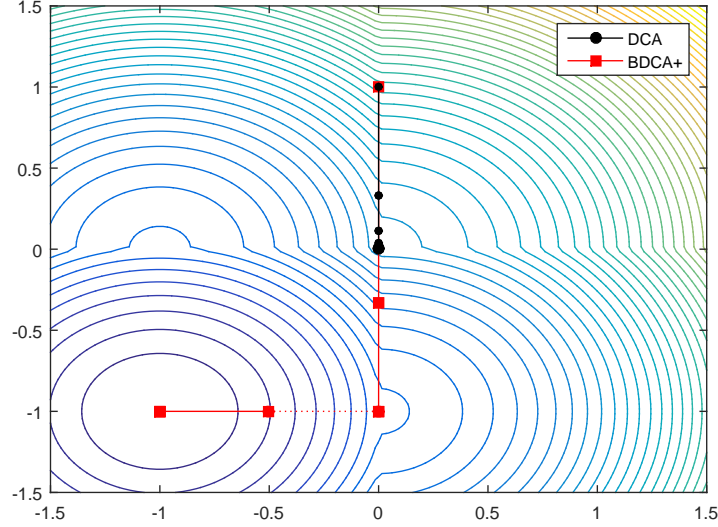


Figure 1: Illustration of Example 3.1.

Consider a corresponding DC decomposition $\phi = g - h$ of ϕ with

$$g(x, y) := \frac{3}{2}(x^2 + y^2) + x + y \quad \text{and} \quad h(x, y) := |x| + |y| + \frac{1}{2}(x^2 + y^2).$$

Observe that g and h satisfy Assumptions 1 to 3. It can be easily checked that ϕ has four critical points, namely $(0, 0)$, $(-1, 0)$, $(0, -1)$ and $(-1, -1)$, of which only the latter is a d-stationary point (and also the global minimum).

In Figure 1 we show the iterations generated by DCA (Algorithm 1) and BDCA+ (Algorithm 3) from the same starting point $x_0 = (0, 1)$. The DCA converges to the critical point $(0, 0)$. The BDCA escapes from this point but still gets stuck at $(0, -1)$, which is also a critical point which is not d-stationary. After applying once the DFO scheme (dashed line), we observe that BDCA successfully converges to the d-stationary point $(-1, -1)$, which is in fact the global minimum of the problem.

To demonstrate the advantage of BDCA+ we compute the number of instances, out of one million random starting points uniformly distributed in $[-1.5, 1.5] \times [-1.5, 1.5]$, that each algorithm has converged to each of the four critical points. The results are summarized in Table 1.

	$(-1, -1)$	$(-1, 0)$	$(0, -1)$	$(0, 0)$
DCA	249,821	250,671	249,944	249,564
BDCA	996,221	1,897	1,882	0
BDCA+	1,000,000	0	0	0

Table 1: For one million random starting points in $[-1.5, 1.5] \times [-1.5, 1.5]$, we count the sequences generated by DCA, BDCA and BDCA+ converging to each of the four d-stationary points

From Table 1, we observe that DCA converged to each of the four critical points with the same probability, while BDCA converged to the global minimum in 99.6% of the instances. The best results were obtained by BDCA+, which always converged to the global minimum $(-1, -1)$.

4 Numerical experiments

In this section, we provide the results of some numerical tests to compare the performance of BDCA+ (Algorithm 3) and the classical DCA (Algorithm 1). To this aim we turn to the same challenging clustering problem tested in [3, Section 5.1], where both algorithms have troubles for finding good solutions due to an abundance of critical points. A different algorithm based on the DC programming approach to solve this problem was introduced in [4]. This algorithm is also proved to converge to d-stationary points.

All the codes were written in Python 2.7 and the tests were run on a desktop of Intel Core i7-4770 CPU 3.40GHz with 32GB RAM, under Windows 10 (64-bit). The following strategies have been followed in all the experiments:

- The trial step size $\bar{\lambda}_k$ in the boosting step of BDCA (line 6 in Algorithms 2 and 3) was chosen to be self-adaptive, as in [3, Section 5], which proceeds as follows:
 1. Set $\bar{\lambda}_0 = 0$ and fix any $\gamma > 1$.
 2. Choose any $\bar{\lambda}_1 > 0$ and obtain λ_1 by backtracking.
 3. For $k \geq 2$,

if $(\lambda_{k-2} = \bar{\lambda}_{k-2} \text{ and } \lambda_{k-1} = \bar{\lambda}_{k-1})$ **then**
 set $\bar{\lambda}_k := \gamma \lambda_{k-1}$;
else set $\bar{\lambda}_k := \lambda_{k-1}$;
 and obtain λ_k by backtracking.

- In our numerical tests we observed that the accepted step sizes μ in the DFO step of Algorithm 3 usually decrease (nearly always). For this reason, we used $\eta := \frac{1}{\beta_2}$ and $\tau := \varepsilon_2$ in the choice of the initial value of μ at line 12 in Algorithm 3. By this way, we allow a slight increase in the value of the step size with respect to the previous one, while we can avoid wasting too much time in this backtracking.
- We tested the three positive basis presented in Example 2.1. Surprisingly, the basis with equally spaced angles D_3 in (3c) performs worse than the others in our test problem. In fact, the best choice was the basis D_1 in (3a), and this is the one we have employed in all the experiments throughout this section.
- We used the parameter setting as $\alpha := 0.0001$, $\varepsilon_1 := 10^{-8}$, $\varepsilon_2 := 10^{-4}$, $\bar{\mu} := 10$, $\gamma = 2$, $\bar{\lambda}_1 := 10$, $\beta_1 := 0.25$ and $\beta_2 := 0.5$.

The Minimum Sum-of-Squares Clustering Problem: Given a collection of n points, $\{a^1, a^2, \dots, a^n \in \mathbb{R}^m\}$, the goal of *clustering* is to group them in k disjoint sets (called *clusters*), $\{A^1, A^2, \dots, A^k\}$, under an optimal criterion. For each cluster A_j , $j = 1, 2, \dots, k$, consider its centroid x^j as a representative. The *Minimum Sum-of-Squares Clustering* criterion asks for the configuration that minimizes the sum of squared distances of each point to its closest centroid, i.e. the solution to the optimization problem

$$\min_{x^1, \dots, x^k \in \mathbb{R}^m} \left\{ \varphi(x^1, \dots, x^k) := \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, k} \|x^j - a^i\|^2 \right\}. \quad (5)$$

We can rewrite the objective in (5) as a DC function (see [4, 8, 11]) with

$$g(x^1, \dots, x^k) := \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k \|x^j - a^i\|^2 + \frac{\rho}{2} \sum_{j=1}^k \|x^j\|^2,$$

$$h(x^1, \dots, x^k) := \frac{1}{n} \sum_{i=1}^n \max_{j=1, \dots, k} \sum_{t=1, t \neq j}^k \|x^t - a^i\|^2 + \frac{\rho}{2} \sum_{j=1}^k \|x^j\|^2;$$

where g and h satisfy Assumptions 1 to 3 for all $\rho > 0$ (in our tests, we took $\rho = \frac{1}{nk}$).

Data Set and Experiments: Our data set is the same one considered in [3], which consists of the location of 4001 Spanish cities in the peninsula with more than 500 inhabitants¹. In Figure 2 we compare the iterations generated by DCA and BDCA+ for finding a partition into 20 clusters from the same random starting point $x_0 \in \mathbb{R}^{2 \times 20}$ (marked with a black cross). We observe that DCA converges to a critical point which is far from being optimal, as there are three clusters without any cities assigned. On the other hand, although BDCA apparently converges to the same critical point, the DFO step allows BDCA+ to escape from points which are not d-stationary and reach a better solution.

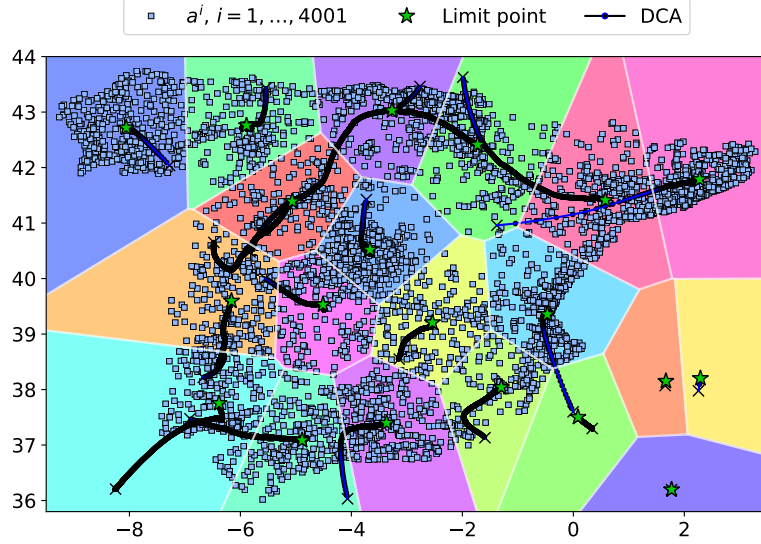
To corroborate these results, we repeated the experiment for different number of clusters $k \in \{20, 40, 60, 80\}$. For each of these values, we run DCA and BDCA+ from 50 random starting points. The results are shown in Figure 3, where we can clearly observe that BDCA+ outperforms DCA, not only in terms of the objective value attained, but even in running time. Observe that it is not really fair to compare the running time of DCA and BDCA+, because DCA simply stops at a critical point without incorporating the time-consuming DFO step that guarantees d-stationarity. Nonetheless, the speedup obtained by the line search of BDCA allows BDCA+ to still converge faster than DCA in most of the instances. As expected, BDCA+ becomes slower as the size of the problem increases, due to the DFO step. Despite that, notice that for 80 clusters the best solution provided by DCA among the 50 instances is still worse than the worst solution obtained by BDCA+. That is, any of the runs of BDCA+ was able to obtain a better solution than 50 restarts of DCA.

5 Concluding remarks

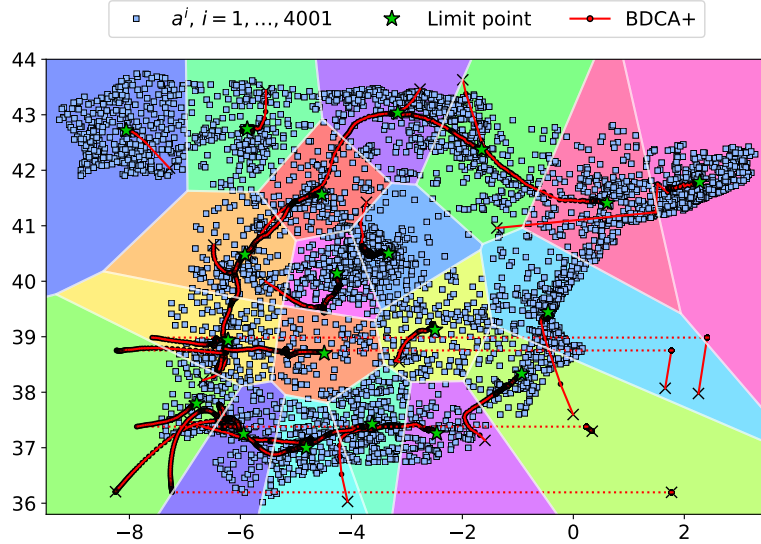
We have proposed a combination between the Boosted DC Algorithm (BDCA) and a simple direct search Derivative-Free Optimization (DFO) technique for minimizing the difference of two convex functions, the first of which is assumed to be smooth. The BDCA is used for minimizing the objective function, while the DFO step permits to force the iteration to converge to d-stationary points (i.e. to points where there exists no descent direction), rather than just critical points.

The good behavior of the new algorithm, called BDCA+, has been demonstrated by numerical experiments in a clustering problem. The new scheme generates better solutions than the classical DCA in nearly all the instances tested. Moreover, this improvement in the quality of the solutions has not caused an important loss in the

¹The data can be retrieved from the Spanish National Center of Geographic Information at <http://centrodedescargas.cnig.es>.



(a) DCA, objective value obtained: 0.4778



(b) BDCA+, objective value obtained: 0.4076

Figure 2: Iterations and limit points generated by DCA and BDCA+ for grouping the Spanish cities in the peninsula into 20 clusters from the same random starting point. The DFO step in line 14 of Algorithm 3 was run 10 times (these steps are marked with a dashed line).

time spent by the algorithm. In fact, BDCA+ was faster than DCA in most of the cases, thanks to the large acceleration achieved by the line search boosting step of BDCA.

Acknowledgements We thank the referees for their constructive comments which helped us improve the presentation of the paper. The first author was supported by MINECO of Spain and ERDF of EU, as part of the Ramón y Cajal program (RYC-2013-13327) and the grants MTM2014-59179-C2-1-P and PGC2018-097960-B-C22. The second author was supported by MINECO of Spain and ESF of EU under the program “Ayudas para contratos predoctorales para la formación de doctores 2015” (BES-2015-073360). The last author was supported by FWF (Austrian Science Fund), project

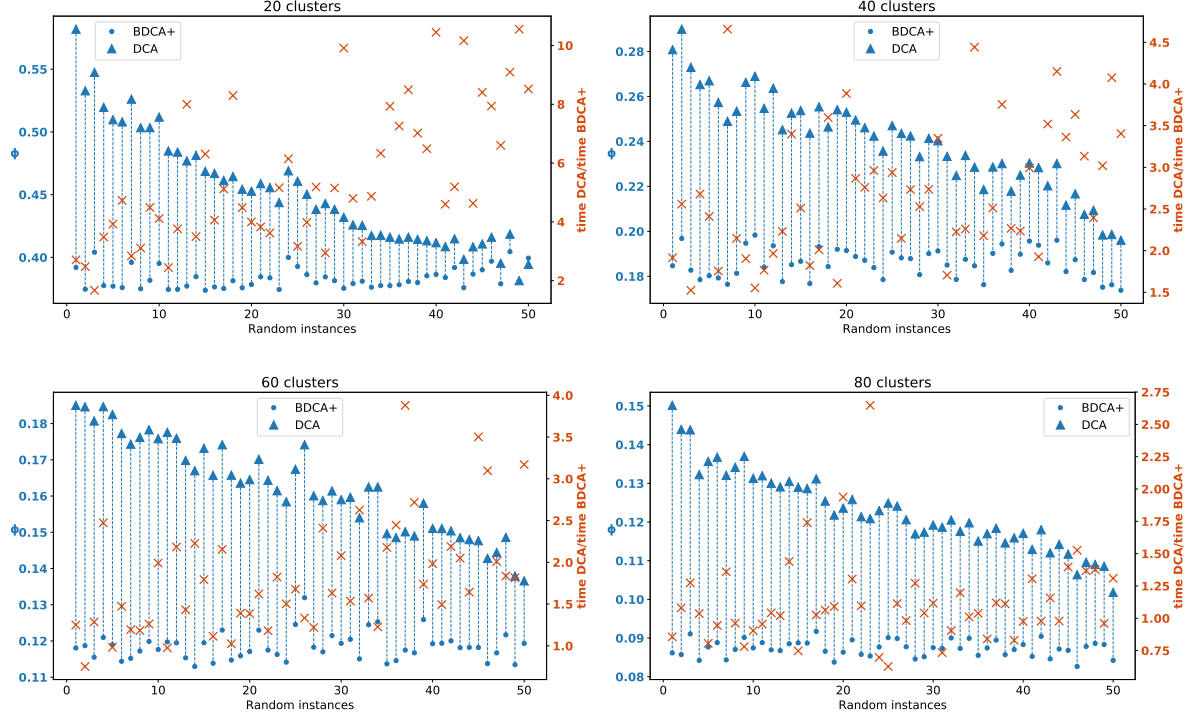


Figure 3: Comparison between the DCA and the BDCA+ for classifying the Spanish cities in the peninsula into k clusters for $k \in \{20, 40, 60, 80\}$. For each of these values, both algorithms were run from 50 random starting points. We represent the objective value achieved in the limit point by each algorithm (left axis, in blue), as well as the ratio between the CPU time required by DCA with respect to the one needed by BDCA+ (right axis, orange crosses). Instances were sorted on the x-axis in descending order according to the gap between the objective values at the limit points found by the algorithms.

M2499-N32 and Vietnam National Foundation for Science and Technology Development (NAFOSTED) project 101.01-2019.320.

References

- [1] Aragón, F.J., Goberna, M.A., López, M.A., Rodríguez, M.M.L.: *Nonlinear Optimization*, Springer Undergraduate Texts in Mathematics and Technology (2019)
- [2] Aragón Artacho, F.J., Fleming, R., Vuong, P.T.: Accelerating the DC algorithm for smooth functions. *Math. Program.* 169(1), 95–118 (2018)
- [3] Aragón Artacho, F.J., Vuong, P.T.: The boosted DC algorithm for nonsmooth functions. To appear in *SIAM J. Optim.* *ArXiv preprint: 1812.06070* (2019)
- [4] Bagirov, A.M., Taheri, S., Ugon, J.: Nonsmooth DC programming approach to the minimum sum-of-squares clustering problems. *Pattern Recogn.* 53, 12–24 (2016)
- [5] Gaudioso, M., Giallombardo, G., Miglionico, G., Bagirov, A.M.: Minimizing nonsmooth DC functions via successive DC piecewise-affine approximations. *J. Glob. Optim.* 71(1), 37–55 (2018)

- [6] Beck, A., Hallak, N.: On the convergence to stationary points of deterministic and randomized feasible descent directions methods. *SIAM J. Optim.* 30(1), 56–79 (2020)
- [7] Conn, A.R., Scheinberg, K., Vicente, L.N.: *Introduction to derivative-free optimization*, MPS-SIAM Series on Optimization vol. 8 (2009)
- [8] Cuong, T.H., Yao, J.C., Yen, N.D.: Qualitative properties of the Minimum Sum-of-Squares Clustering problem. *ArXiv preprint: 1810.02057* (2018)
- [9] Le Thi, H.A., Pham Dinh, T.: DC Programming and DCA: Thirty years of developments. *Math. Program.* 169(1), 5–68 (2018)
- [10] Le Thi, H.A., Pham Dinh, T.: The DC (difference of convex functions) programming and DCA revisited with DC models of real world nonconvex optimization problems. *Ann. Oper. Res.* 133(1-4), 23–46 (2005)
- [11] Ordin, B., Bagirov, A.M.: A heuristic algorithm for solving the minimum sum-of-squares clustering problems. *J. Glob. Optim.* 61, 341–361 (2015)
- [12] de Oliveira, W.: Proximal bundle methods for nonsmooth DC programming. *J. Glob. Optim.* 75, 523–563 (2019)
- [13] de Oliveira, W., Tcheou, M.P.: An inertial algorithm for DC programming. *Set-Valued Var. Anal.* 27, 895–919(2019)
- [14] Pham Dinh, T., Le Thi, H.A.: Convex analysis approach to DC programming: theory, algorithms and applications. *Acta Mathematica Vietnamica*, 22, 289–355 (1997)
- [15] Pham Dinh T., Souad, E.B.: Algorithms for solving a class of nonconvex optimization problems. Methods of subgradients. In J.-B. Hiriart-Urruty, editor, *FERMAT Days 85: Mathematics for Optimization, volume 129 of North-Holland Mathematics Studies*, pp. 249–271. Elsevier (1986)
- [16] Toland, J.F.: On subdifferential calculus and duality in non-convex optimization. *Bull. Soc. Math. Fr. Mém.* 60 (Proc. Colloq., Pau 1977), 177–183 (1979)
- [17] Rockafellar, R.T.: *Convex Analysis*, Princeton University Press (1972)