

Optimization

A Journal of Mathematical Programming and Operations Research

ISSN: 0233-1934 (Print) 1029-4945 (Online) Journal homepage: <https://www.tandfonline.com/loi/gopt20>

Clustering and multifacility location with constraints via distance function penalty methods and dc programming

Nguyen Mau Nam, Nguyen Thai An, Sam Reynolds & Tuyen Tran

To cite this article: Nguyen Mau Nam, Nguyen Thai An, Sam Reynolds & Tuyen Tran (2018) Clustering and multifacility location with constraints via distance function penalty methods and dc programming, Optimization, 67:11, 1869-1894, DOI: [10.1080/02331934.2018.1510498](https://doi.org/10.1080/02331934.2018.1510498)

To link to this article: <https://doi.org/10.1080/02331934.2018.1510498>



Published online: 02 Sep 2018.



Submit your article to this journal [↗](#)



Article views: 62



View Crossmark data [↗](#)



Clustering and multifacility location with constraints via distance function penalty methods and dc programming

Nguyen Mau Nam^a, Nguyen Thai An^b, Sam Reynolds^a and Tuyen Tran^a

^aFariborz Maseeh Department of Mathematics and Statistics, Portland State University, Portland, OR, USA; ^bInstitute of Research and Development, Duy Tan University, Danang, Vietnam

ABSTRACT

This paper is a continuation of our effort in using mathematical optimization involving DC programming in clustering and multifacility location. We study a penalty method based on distance functions and apply it particularly to a number of problems in clustering and multifacility location in which the centers to be found must lie in some given set constraints. We also provide numerical examples to test our method.

ARTICLE HISTORY

Received 21 April 2018
Accepted 22 July 2018

KEYWORDS

Clustering; DC programming; Nesterov's smoothing techniques; k-mean algorithm

AMS SUBJECT CLASSIFICATIONS

49J52; 49J53; 90C31

1. Introduction

In the current time of “big data”, clustering is a very important problem that helps classify data in many fields such as machine learning, pattern recognition, image analysis, data compression, and computer graphics. Given a finite number of data points with a measurement distance, a centroid-based clustering problem seeks a finite number of cluster centers with each data point assigned to the nearest cluster center in a way that a certain measurement distance is minimized.

It is well-known that the k —means algorithm is one of the simplest clustering algorithms, providing an easy way to classify a given data set through a certain number of clusters. However, it possesses certain drawbacks: the k —mean algorithm depends heavily on the initial choice of cluster centers; there is no guarantee that the k —means algorithm converges to a global optimal solution; the number of clusters k is an input parameter: an inappropriate choice of k may yield poor results; the results depend heavily on the measurement distance; the algorithm may not be applicable for handling effectively constraints imposed on the cluster centers.

In our recent research, we further the pioneering works by Pham Dinh Tao, Le Thi Hoai An and others from [1,2] in using the mathematical programming approach for clustering, aiming at providing an alternative to the k —means

algorithm and coping with its drawbacks; see [1,3,4]. The mathematical programming approach is very promising as optimization techniques for minimizing nonconvex optimization problems have been of great interest with significant progress over the past few years. In addition, it is possible to use derivative-free methods for initializations in the DCA and enhance the effectiveness of gradient/subgradient-based nonconvex algorithms. Our method using Nesterov's smoothing techniques and the DCA, an algorithm for minimizing differences of convex functions, allows us to solve clustering and multifacility location problems in many different settings involving different norms, bilevel clustering, and set clustering.

The main focus of this paper is on solving a number of clustering and multifacility location problems with constraints. We use a penalty method with squared Euclidean distance functions to convert constrained problems to unconstrained problems. Then appropriate DC decompositions and the DCA are used to minimize the penalized objective functions. In the case where the measurement distance is defined by the Euclidean norm instead of the squared Euclidean norm, we use Nesterov's smoothing techniques for *reducing the nonsmoothness* of the model and for providing a DC decomposition that is favorable for applying the DCA. Our method opens up the possibility of using distance function penalty methods for other problems of DC programming.

The paper is organized as follows. In Section 2, we present basic tools of convex analysis and optimization used throughout the paper. The analysis of a penalty method based on squared distance functions is presented in Section 3. Section 4 is devoted to solving clustering problems with constraints in which the measurement distance is defined by the squared Euclidean norm. In Section 5, we study a new model of clustering with constraints that involves sets. In Section 6, we study clustering problems with constrained and the measurement distance defined by the Euclidean norm. These problems belong to the class of continuous multifacility location problems with constraints. Finally, numerical examples are presented in Section 7.

Throughout the paper, we use $\langle \cdot, \cdot \rangle$ to denote the inner product and use $\| \cdot \|$ to denote the associated Euclidean norm in \mathbb{R}^d . For the subset Ω of \mathbb{R}^d , the set $\text{conv}(\Omega)$ is the convex hull of Ω , i.e. the smallest convex set in \mathbb{R}^d that contains Ω .

2. Preliminaries

In this section, we present basic tools of analysis and optimization used in the sequel. The readers are referred to [1,5–7] for more details.

Let $f: \mathbb{R}^d \rightarrow (-\infty, \infty]$ be a convex function. An element $v \in \mathbb{R}^d$ is called a *subgradient* of f at $\bar{x} \in \text{dom}(f) = \{x \in \mathbb{R}^d | f(x) < \infty\}$ if it satisfies

$$\langle v, x - \bar{x} \rangle \leq f(x) - f(\bar{x}) \quad \text{for all } x \in \mathbb{R}^d.$$

The set of all such elements v is called the *subdifferential* of f at \bar{x} and is denoted by $\partial f(\bar{x})$. If $\bar{x} \notin \text{dom}(f)$, we set $\partial f(\bar{x}) = \emptyset$. This subdifferential concept possesses many calculus rules that are important for applications. In particular, for a finite number of convex functions $f_i: \mathbb{R}^d \rightarrow (-\infty, \infty]$, $i = 1, \dots, m$, we have the following sum rule:

$$\partial(f_1 + \dots + f_m)(\bar{x}) = \partial f_1(\bar{x}) + \dots + \partial f_m(\bar{x}) \quad \text{for all } \bar{x} \in \mathbb{R}^d$$

provided that $\bigcap_{i=1}^m \text{ri}(\text{dom}(f_i)) \neq \emptyset$. Here $\text{ri}(\Omega)$ stands for the *relative interior* of Ω ; see, e.g. [6, Definition 1.68].

If $f = \max_{i=1, \dots, m} f_i$, and f_i is continuous at \bar{x} for every $i = 1, \dots, m$, then for any $\bar{x} \in \mathbb{R}^d$ we have the following maximal rule:

$$\partial f(\bar{x}) = \text{conv}\left(\bigcup_{i \in I(\bar{x})} \partial f_i(\bar{x})\right), \quad (1)$$

where $I(\bar{x}) = \{i | f_i(\bar{x}) = f(\bar{x})\}$.

Given a nonempty closed convex subset Ω of \mathbb{R}^d with $\bar{x} \in \Omega$, the *normal cone* to Ω at \bar{x} is defined by

$$N(\bar{x}; \Omega) = \{v \in \mathbb{R}^d \mid \langle v, x - \bar{x} \rangle \leq 0 \text{ for all } x \in \Omega\}.$$

If $\bar{x} \notin \Omega$, we set $N(\bar{x}, \Omega) = \emptyset$. It is well-known that an element $\bar{x} \in \mathbb{R}^d$ is an absolute minimizer of a convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ on Ω if and only if \bar{x} is a local minimizer of f on Ω . Moreover, this happens if and only if the following optimality condition holds: $0 \in \partial f(\bar{x}) + N(\bar{x}; \Omega)$.

Let $\Theta \subset \mathbb{R}^d$ be a nonempty set (not necessarily convex). The *distance function* to Θ is defined by

$$d(x; \Theta) = \inf \{\|x - w\| \mid w \in \Theta\}, \quad x \in \mathbb{R}^d.$$

The *Euclidean project* from $x \in \mathbb{R}^d$ to Θ is the set

$$P(x; \Theta) = \{w \in \Theta \mid d(x; \Theta) = \|x - w\|\}.$$

We can show that if Θ is a nonempty closed set, then $P(x; \Theta)$ is nonempty, and it is a singleton if we assume in addition that Θ is convex. We can also show that if Θ is a convex set and $w \in P(x; \Theta)$, then $x - w \in N(w; \Theta)$.

Another tool we will use in the paper is the notion of *Fenchel conjugates*. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a function. The Fenchel conjugate of f is defined by

$$f^*(y) = \sup \{\langle y, x \rangle - f(x) \mid x \in \mathbb{R}^d\}, \quad y \in \mathbb{R}^d.$$

Note that $f^*: \mathbb{R}^d \rightarrow (-\infty, \infty]$ is an extended-real-valued convex function. Suppose further that f is convex, then the *Fenchel-Moreau theorem* states that

$(f^*)^* = f$. Based on this theorem, we have the following relation between the subgradients of f and its Fenchel conjugate:

$$x \in \partial f^*(y) \iff y \in \partial f(x). \quad (2)$$

The notions of subgradients and Fenchel conjugates provide mathematical foundation for the DCA introduced below. Given a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ with the *DC decomposition* $f = g - h$, where $g, h: \mathbb{R}^d \rightarrow \mathbb{R}$ are convex functions, the DCA introduced by Pham Dinh Tao described in what follows is a simple but effective algorithm for minimizing the function f ; see [8,9].

Algorithm 1 The DCA

INPUT: $x_1, N \in \mathbb{N}$

for $p = 1, \dots, N$ **do**

 Find $y_p \in \partial h(x_p)$

 Find $x_{p+1} \in \partial g^*(y_p)$

end for

OUTPUT: x_{N+1}

For convenience, define the *data matrix* $\mathbf{A} \in \mathbb{R}^{m \times d}$ as the matrix whose i^{th} row is $a^i \in \mathbb{R}^d$ for $i = 1, \dots, m$. Similarly, we define the *variable matrix* $\mathbf{X} \in \mathbb{R}^{k \times d}$ as the matrix whose ℓ^{th} row is $x^\ell \in \mathbb{R}^d$ for $\ell = 1, \dots, k$. We equip the linear space $\mathbb{R}^{k \times d}$ with the inner product $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{trace}(\mathbf{X}^T \mathbf{Y})$. Recall that the *Frobenius norm* on $\mathbb{R}^{k \times d}$ is defined by

$$\|\mathbf{X}\|_F = \langle \mathbf{X}, \mathbf{X} \rangle^{1/2} = \left(\sum_{\ell=1}^k \langle x^\ell, x^\ell \rangle \right)^{1/2} = \left(\sum_{\ell=1}^k \|x^\ell\|^2 \right)^{1/2}.$$

Observe that the square of the Frobenius norm is differentiable with

$$\nabla \|\mathbf{X}\|_F^2 = 2\mathbf{X} \quad \text{for } \mathbf{X} \in \mathbb{R}^{k \times d}.$$

Let $\Omega^\ell \subset \mathbb{R}^d$ for $\ell = 1, \dots, k$ be nonempty closed convex sets and let $\mathbf{\Omega} = \Omega^1 \times \Omega^2 \times \dots \times \Omega^k$. For $\mathbf{X} \in \mathbb{R}^{k \times d}$, the projection from \mathbf{X} to $\mathbf{\Omega}$ is the matrix \mathbf{Y} whose ℓ^{th} row is $y^\ell = P(x^\ell; \Omega^\ell)$. We thus have

$$[d(\mathbf{X}; \mathbf{\Omega})]^2 = \|\mathbf{X} - \mathbf{Y}\|_F^2 = \sum_{\ell=1}^k \|x^\ell - y^\ell\|^2 = \sum_{\ell=1}^k d(x^\ell; \Omega^\ell)^2. \quad (3)$$

3. A penalty method via distance functions

In this section, we study a penalty method using distance functions for solving constrained optimization problems and apply them specifically to DC programming. This method is based on the *quadratic penalty method*; see [10,11]. Let

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a function and let Ω_i for $i = 1, \dots, q$ be nonempty closed subsets of \mathbb{R}^d with $\bigcap_{i=1}^q \Omega_i \neq \emptyset$. Consider the optimization problem:

$$\begin{aligned} \min \quad & f(x) \\ \text{subject to} \quad & x \in \bigcap_{i=1}^q \Omega_i. \end{aligned} \quad (4)$$

Let us first study the relation between this problem and the unconstrained problem given by

$$\min f_\lambda(x) = f(x) + \frac{\lambda}{2} \sum_{i=1}^q [d(x; \Omega_i)]^2, \quad x \in \mathbb{R}^d. \quad (5)$$

The theorem below provides a relation between optimal solutions of the constrained optimization problem (4) and the unconstrained optimization problem (5) obtained by a penalty method based on distance functions. The proof follows [11, Theorem 17.1].

Theorem 3.1: Consider (4) in which $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a l.s.c. function. Suppose that (4) has an optimal solution. If $\lim_{n \rightarrow \infty} \lambda_n = \infty$ and $x_n \in \mathbb{R}^d$ is an absolute minimizer of the function f_{λ_n} defined in (5) for all $n \in \mathbb{N}$, then every subsequential limit of $\{x_n\}$ is a solution of (4).

Proof: Let $\bar{x} \in \mathbb{R}^d$ be an optimal solution of (4). That means $\bar{x} \in \Omega_i$ for $i = 1, \dots, q$ and

$$f(\bar{x}) \leq f(x) \quad \text{whenever } x \in \Omega_i \quad \text{for all } i = 1, \dots, q.$$

Since $x_n \in \mathbb{R}^d$ is an absolute minimizer of the function f_{λ_n} ,

$$f_{\lambda_n}(x_n) \leq f_{\lambda_n}(\bar{x}).$$

This implies, with the observation that $d(\bar{x}; \Omega_i) = 0$ for $i = 1, \dots, q$, that

$$f(x_n) + \frac{\lambda_n}{2} \sum_{i=1}^q [d(x_n; \Omega_i)]^2 \leq f(\bar{x}). \quad (6)$$

Then

$$\sum_{i=1}^q [d(x_n; \Omega_i)]^2 \leq \frac{2}{\lambda_n} (f(\bar{x}) - f(x_n)) \quad \text{for all } n \in \mathbb{N}.$$

Let $x^* \in \mathbb{R}^d$ be a subsequential limit of $\{x_n\}$. Without loss of generality, we can assume that $\lim_{n \rightarrow \infty} x_n = x^*$. By the continuity of the distance function and the

lower semicontinuity of f ,

$$\sum_{i=1}^q [d(x^*; \Omega_i)]^2 = \lim_{n \rightarrow \infty} \sum_{i=1}^q [d(x_n; \Omega_i)]^2 \leq \liminf_{n \rightarrow \infty} \frac{2}{\lambda_n} (f(\bar{x}) - f(x_n)) \leq 0.$$

It follows that $d(x^*; \Omega_i) = 0$, and so $x^* \in \Omega_i$ for $i = 1, \dots, q$. In addition, by (6) and the lower semicontinuity of f we have

$$f(x^*) \leq \liminf_{n \rightarrow \infty} f(x_n) \leq \liminf_{n \rightarrow \infty} \left(f(x_n) + \frac{\lambda_n}{2} \sum_{i=1}^q [d(x_n; \Omega_i)]^2 \right) \leq f(\bar{x}).$$

Therefore, x^* is an optimal solution of (4). ■

Now we discuss a direct consequence of Theorem 3.1 that will be used in the sequel. Let $F: \mathbb{R}^{k \times d} \rightarrow \mathbb{R}$ be a function and let Ω_i^ℓ for $\ell = 1, \dots, k$ and $i = 1, \dots, q$ be nonempty closed subsets of \mathbb{R}^d . Consider the problem

$$\begin{aligned} \min \quad & F(x^1, \dots, x^k) \\ \text{subject to} \quad & x^\ell \in \bigcap_{i=1}^q \Omega_i^\ell \quad \text{for } \ell = 1, \dots, k. \end{aligned} \quad (7)$$

We now clarify the relation between this problem and the unconstrained problem given by

$$\begin{aligned} \min \quad & F_\lambda(x^1, \dots, x^k) = F(x^1, \dots, x^k) + \frac{\lambda}{2} \sum_{\ell=1}^k \sum_{i=1}^q [d(x^\ell; \Omega_i^\ell)]^2 \\ & x^\ell \in \mathbb{R}^d \quad \text{for } \ell = 1, \dots, k. \end{aligned} \quad (8)$$

In what follows, we identify $\mathbf{X} = (x^1, \dots, x^k) \in \mathbb{R}^{k \times d}$ with the matrix $\mathbf{X} \in \mathbb{R}^{k \times d}$, whose ℓ^{th} row is x^ℓ for $\ell = 1, \dots, k$.

Corollary 3.2: Consider (7) in which $F: \mathbb{R}^{k \times d} \rightarrow \mathbb{R}$ is a l.s.c. function. Suppose that (7) has an optimal solution. If $\lim_{n \rightarrow \infty} \lambda_n = \infty$ and $X_n = (x_n^1, \dots, x_n^k) \in \mathbb{R}^{k \times d}$ is an absolute minimizer of the function F_{λ_n} , then every subsequential limit of $\{X_n\}$ is a solution of (7).

Proof: Let $\mathbf{X} = (x^1, \dots, x^k) \in \mathbb{R}^{k \times d}$ and let $\Omega_i = \Omega_i^1 \times \dots \times \Omega_i^k \subset \mathbb{R}^{k \times d}$ for $i = 1, \dots, q$. Note that

$$\bigcap_{i=1}^q \Omega_i = \bigcap_{i=1}^q \prod_{\ell=1}^k \Omega_i^\ell = \prod_{\ell=1}^k \bigcap_{i=1}^q \Omega_i^\ell.$$

It follows that $x^\ell \in \bigcap_{i=1}^q \Omega_i^\ell$ for $\ell = 1, \dots, k$ if and only if $\mathbf{X} \in \bigcap_{i=1}^q \Omega_i$, and thus (7) reduces to the following optimization problem:

$$\begin{aligned} \min \quad & F(\mathbf{X}) \\ \text{subject to} \quad & \mathbf{X} \in \bigcap_{i=1}^q \Omega_i. \end{aligned}$$

Based on (3), we can rewrite the objective function F_λ in (8) as follows

$$F_\lambda(\mathbf{X}) = F(\mathbf{X}) + \frac{\lambda}{2} \sum_{i=1}^q [d(\mathbf{X}; \Omega_i)]^2.$$

The conclusion now follows directly from Theorem 3.1. ■

Let us continue with a known result on DC decompositions of squared distance functions. The proof of the following result can be found in [4, Proposition 5.1].

Proposition 3.3: *Let Ω be a nonempty closed set in \mathbb{R}^d (not necessarily convex). Define the function*

$$\varphi_\Omega(x) = \sup \{ \langle 2x, w \rangle - \|w\|^2 \mid w \in \Omega \} = 2 \sup \{ \langle x, w \rangle - \frac{1}{2} \|w\|^2 \mid w \in \Omega \}.$$

Then we have the following conclusions:

- (i) *The function φ_Ω is always convex. If we assume in addition that Ω is convex, then φ_Ω is differentiable with $\nabla \varphi_\Omega(x) = 2P(x; \Omega)$.*
- (ii) *The function $f(x) = [d(x; \Omega)]^2$ is a DC function with $f(x) = \|x\|^2 - \varphi_\Omega(x)$ for all $x \in \mathbb{R}^d$.*

We now consider (4) in which $f(x) = g(x) - h(x)$ is a DC function where $g, h: \mathbb{R}^d \rightarrow \mathbb{R}$ are convex functions. We also assume additionally that all constraint sets are convex and satisfy $\bigcap_{i=1}^q \text{ri}(\Omega_i) \neq \emptyset$. By [12, Theorem 5.3], this condition ensures that

$$N(\bar{x}; \bigcap_{i=1}^q \Omega_i) = \sum_{i=1}^q N(\bar{x}; \Omega_i) \quad \text{for every } \bar{x} \in \bigcap_{i=1}^q \Omega_i. \quad (9)$$

Recall from [9] that an element $\bar{x} \in \mathbb{R}^d$ is a *critical point* of a DC function f with DC decomposition $f = g - h$ if $\partial g(\bar{x}) \cap \partial h(\bar{x}) \neq \emptyset$. Observe that (4) can be written as an unconstrained optimization problem using the indicator function as follows:

$$\min \left(g(x) + \delta(x, \bigcap_{i=1}^q \Omega_i) \right) - h(x), \quad x \in \mathbb{R}^d.$$

Define $v(x) = g(x) + \delta(x, \bigcap_{i=1}^q \Omega_i)$ for $x \in \mathbb{R}^d$. By (9),

$$\partial v(\bar{x}) = \partial g(\bar{x}) + N(\bar{x}; \bigcap_{i=1}^q \Omega_i) = \partial g(\bar{x}) + \sum_{i=1}^q N(\bar{x}; \Omega_i) \quad \text{for all } \bar{x} \in \mathbb{R}^d.$$

Thus, we call an element $\bar{x} \in \mathbb{R}^d$ a critical point of (4) if

$$\left(\partial g(\bar{x}) + \sum_{i=1}^q N(\bar{x}; \Omega_i) \right) \cap \partial h(\bar{x}) \neq \emptyset. \quad (10)$$

The objective function of (5) now becomes

$$f_\lambda(x) = g(x) + \frac{\lambda}{2} \sum_{i=1}^q d(x; \Omega_i)^2 - h(x).$$

Using Proposition 3.3, we have

$$f_\lambda(x) = \left(g(x) + \frac{\lambda q}{2} \|x\|^2 \right) - \left(h(x) + \frac{\lambda}{2} \sum_{i=1}^q \varphi_{\Omega_i}(x) \right) = \tilde{g}_\lambda(x) - \tilde{h}_\lambda(x), \quad (11)$$

where \tilde{g}_λ and \tilde{h}_λ are functions defined on \mathbb{R}^d by

$$\tilde{g}_\lambda(x) = g(x) + \frac{\lambda q}{2} \|x\|^2 \text{ and } \tilde{h}_\lambda(x) = h(x) + \frac{\lambda}{2} \sum_{i=1}^q \varphi_{\Omega_i}(x), \quad x \in \mathbb{R}^d.$$

Proposition 3.4: Suppose that $\lim_{n \rightarrow \infty} \lambda_n = \infty$ and x_n is a critical point of the DC function $f_{\lambda_n} = \tilde{g}_{\lambda_n} - \tilde{h}_{\lambda_n}$ given in (11). Then every subsequential limit of the sequence $\{x_n\}$ is a critical point of (4).

Proof: Since x_n is a critical point of f_{λ_n} and by Proposition 3.3, there exist $v_n \in \partial g(x_n)$ and $w_n \in \partial h(x_n)$ such that

$$v_n + \lambda_n q x_n = w_n + \lambda_n \sum_{i=1}^q P(x_n; \Omega_i). \quad (12)$$

Let \bar{x} be a subsequential limit of $\{x_n\}$. Without loss of generality, we can assume that $\{x_n\}$ converges to \bar{x} . Since any finite convex function is locally Lipschitz continuous, we can assume that both g and h are locally Lipschitz continuous around \bar{x} with Lipschitz constant $L > 0$. Then

$$\|v_n\| \leq L \text{ and } \|w_n\| \leq L \text{ for sufficiently large } n. \quad (13)$$

By (12), (13) and the assumption that $\lambda_n \rightarrow \infty$ as $n \rightarrow \infty$,

$$\sum_{i=1}^q \left(x_n - P(x_n; \Omega_i) \right) = \frac{1}{\lambda_n} (w_n - v_n) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Letting $n \rightarrow \infty$ yields $\sum_{i=1}^q (\bar{x} - P(\bar{x}; \Omega_i)) = 0$, due to the continuity of projection operators onto convex sets. Note also that $\lambda_n \sum_{i=1}^q (x_n - P(x_n; \Omega_i)) \in$

$\sum_{i=1}^q N(x_n; \Omega_i)$. This implies

$$w_n - v_n \in \sum_{i=1}^q N(x_n; \Omega_i) = N(x_n; \bigcap_{i=1}^q \Omega_i).$$

By (13), we can assume without loss of generality that $v_n \rightarrow \bar{v}$ and $w_n \rightarrow \bar{w}$ as $n \rightarrow \infty$. Then by passing to a limit, we have

$$\bar{w} - \bar{v} \in N(\bar{x}; \bigcap_{i=1}^q \Omega_i) = \sum_{i=1}^q N(\bar{x}; \Omega_i).$$

Observe also that $\bar{v} \in \partial g(\bar{x})$ and $\bar{w} \in \partial h(\bar{x})$. Therefore, (10) is satisfied and thus \bar{x} is a critical point of (4). ■

We continue by considering (7) in which

$$F(x^1, \dots, x^k) = G(x^1, \dots, x^k) - H(x^1, \dots, x^k)$$

is a DC function, where $G, H: \mathbb{R}^{k \times d} \rightarrow \mathbb{R}$ are convex functions. From the proof of Corollary 3.2, we can rewrite (7) as

$$\begin{aligned} \min \quad & F(\mathbf{X}) = G(\mathbf{X}) - H(\mathbf{X}) \\ \text{subject to} \quad & \mathbf{X} \in \bigcap_{i=1}^q \Omega_i. \end{aligned}$$

Recall that a point $\bar{\mathbf{X}} = (\bar{x}^1, \dots, \bar{x}^k)$ is called a critical point of this problem if

$$\left(\partial G(\bar{\mathbf{X}}) + \sum_{i=1}^q N(\bar{\mathbf{X}}; \Omega_i) \right) \cap \partial H(\bar{\mathbf{X}}) \neq \emptyset,$$

where $N(\bar{\mathbf{X}}; \Omega_i) = N(\bar{x}^1; \Omega_i^1) \times \dots \times N(\bar{x}^k; \Omega_i^k)$.

For $\Omega \subset \mathbb{R}^{k \times d}$, based on Frobenius norm, we define

$$\varphi_{\Omega}(\mathbf{X}) = \|\mathbf{X}\|_F^2 - d(\mathbf{X}; \Omega)^2 = 2 \sup \left\{ \langle \mathbf{X}, \mathbf{Y} \rangle - \frac{\|\mathbf{Y}\|_F^2}{2} \mid \mathbf{Y} \in \Omega \right\}, \quad \mathbf{X} \in \mathbb{R}^{k \times d}.$$

Observe that if $\Omega_i = \Omega_i^1 \times \Omega_i^2 \times \dots \times \Omega_i^k$ and $\mathbf{X} = (x^1, \dots, x^k) \in \mathbb{R}^{k \times d}$, then

$$d(\mathbf{X}; \Omega_i)^2 = \sum_{\ell=1}^k d(x^\ell; \Omega_i^\ell)^2 = \sum_{\ell=1}^k \left(\|x^\ell\|^2 - \varphi_{\Omega_i^\ell}(x^\ell) \right) = \|\mathbf{X}\|_F^2 - \sum_{\ell=1}^k \varphi_{\Omega_i^\ell}(x^\ell).$$

Therefore, $\varphi_{\Omega_i}(\mathbf{X}) = \sum_{\ell=1}^k \varphi_{\Omega_i^\ell}(x^\ell)$.

In this new notation, the function F_λ in (8) can be rewritten as

$$F_\lambda(\mathbf{X}) = \left(G(\mathbf{X}) + \frac{\lambda q}{2} \|\mathbf{X}\|_F^2 \right) - \left(H(\mathbf{X}) + \sum_{i=1}^q \varphi_{\Omega_i}(\mathbf{X}) \right) = G_1(\mathbf{X}) - H_1(\mathbf{X}),$$

where

$$G_1(\mathbf{X}) = G(\mathbf{X}) + \frac{\lambda q}{2} \|\mathbf{X}\|_F^2 \quad \text{and} \quad H_1(\mathbf{X}) = H(\mathbf{X}) + \sum_{i=1}^q \varphi_{\Omega_i}(\mathbf{X}), \quad \mathbf{X} \in \mathbb{R}^{k \times d}.$$

We also recall that $\bar{\mathbf{X}} \in \mathbb{R}^{k \times d}$ is a critical point of (8) if

$$\partial G_1(\bar{\mathbf{X}}) \cap \partial H_1(\bar{\mathbf{X}}) \neq \emptyset.$$

The proof of the following result is similar to that of Proposition 3.4.

Proposition 3.5: *Suppose that $\lim_{n \rightarrow \infty} \lambda_n = \infty$ and $\mathbf{X}_n = (x_n^1, \dots, x_n^k) \in \mathbb{R}^{k \times d}$ is a critical point of the function F_{λ_n} . Then every subsequential limit of $\{\mathbf{X}_n\}$ is a critical point of (7).*

4. Clustering with constraints

In this section, we study problems of *clustering with constraints* in which the measurement distance is defined by the squared Euclidean norm. We seek k centers $x^1, \dots, x^k \in \mathbb{R}^d$ of m data nodes $a^1, \dots, a^m \in \mathbb{R}^d$ and impose the restriction that each $x^\ell \in \bigcap_{i=1}^q \Omega_i^\ell$ for some nonempty closed convex set $\Omega_i^\ell \subset \mathbb{R}^d$ with $\ell = 1, \dots, k$ and $i = 1, \dots, q$. Here, without loss of generality, we assume that the numbers of constraints for each center is equal to each other. The problem we are concerned with is given by

$$\begin{aligned} \min \quad & \psi(x^1, \dots, x^k) = \sum_{i=1}^m \min_{\ell=1, \dots, k} \|x^\ell - a^i\|^2 \\ \text{subject to} \quad & x^\ell \in \bigcap_{j=1}^q \Omega_j^\ell \quad \text{for } \ell = 1, \dots, k. \end{aligned} \quad (14)$$

This problem can be converted to an unconstrained minimization problem:

$$\begin{aligned} \min f(x^1, \dots, x^k) = & \frac{1}{2} \sum_{i=1}^m \min_{\ell=1, \dots, k} \|x^\ell - a^i\|^2 \\ & + \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q [d(x^\ell; \Omega_i^\ell)]^2, \quad x^1, \dots, x^k \in \mathbb{R}^d, \end{aligned} \quad (15)$$

where $\tau > 0$ is a penalty parameter.

Recall from Proposition 3.3 that for any nonempty closed convex set Ω in \mathbb{R}^d ,

$$[d(x; \Omega)]^2 = \|x\|^2 - \varphi_\Omega(x),$$

where $\varphi_\Omega(x) = 2 \sup\{\langle x, w \rangle - \frac{1}{2} \|w\|^2 \mid w \in \Omega\}$ is a differentiable function with $\nabla \varphi_\Omega(x) = 2P(x; \Omega)$. Let us use the *minimum-sum principle* for k real numbers

α_ℓ for $\ell = 1, \dots, k$:

$$\min_{\ell=1,\dots,k} \alpha_\ell = \sum_{\ell=1}^k \alpha_\ell - \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k \alpha_\ell$$

to obtain a DC decomposition of f as follows

$$\begin{aligned} f(x^1, \dots, x^k) &= \left(\frac{1}{2} \sum_{i=1}^m \sum_{\ell=1}^k \|x^\ell - a^i\|^2 + \frac{\tau q}{2} \sum_{\ell=1}^k \|x^\ell\|^2 \right) \\ &\quad - \left(\frac{1}{2} \sum_{i=1}^m \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k (\|x^\ell - a^i\|)^2 + \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \varphi_{\Omega_i^\ell}(x^\ell) \right). \end{aligned}$$

We see that $f = g - h$ by defining

$$\begin{aligned} g_1(x^1, \dots, x^k) &= \frac{1}{2} \sum_{i=1}^m \sum_{\ell=1}^k \|x^\ell - a^i\|^2, \quad g_2(x^1, \dots, x^k) = \frac{\tau q}{2} \sum_{\ell=1}^k \|x^\ell\|^2, \\ h_1(x^1, \dots, x^k) &= \frac{1}{2} \sum_{i=1}^m \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k \|x^\ell - a^i\|^2, \quad h_2(x^1, \dots, x^k) \\ &= \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \varphi_{\Omega_i^\ell}(x^\ell), \end{aligned}$$

and setting $g = g_1 + g_2$ and $h = h_1 + h_2$.

As discussed in the introduction, we may collect x^j into the variable matrix \mathbf{X} and denote $\Omega_i = \Omega_i^1 \times \Omega_i^2 \times \dots \times \Omega_i^k \in \mathbb{R}^{k \times d}$ for $i = 1, \dots, q$. Then (14) becomes

$$\min \psi(\mathbf{X}) \quad \text{subject to } \mathbf{X} \in \bigcap_{i=1}^q \Omega_i.$$

We also collect a^i into the data matrix \mathbf{A} , and upon doing so we may express g in terms of the Frobenius norm, namely,

$$\begin{aligned} g_1(\mathbf{X}) &= \frac{1}{2} \sum_{i=1}^m \sum_{\ell=1}^k \|x^\ell - a^i\|^2 = \frac{1}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left(\|x^\ell\|^2 - 2\langle x^\ell, a^i \rangle + \|a^i\|^2 \right) \\ &= \frac{m}{2} \sum_{\ell=1}^k \|x^\ell\|^2 - \sum_{i=1}^m \sum_{\ell=1}^k \langle x^\ell, a^i \rangle + \frac{k}{2} \sum_{i=1}^m \|a^i\|^2 \\ &= \frac{m}{2} \|\mathbf{X}\|_F^2 - \langle \mathbf{X}, \mathbf{EA} \rangle + \frac{k}{2} \|\mathbf{A}\|_F^2, \end{aligned}$$

where $\mathbf{E} \in \mathbb{R}^{k \times m}$ is the matrix of ones. In this form, it is easily seen that

$$\nabla g_1(\mathbf{X}) = m\mathbf{X} - \mathbf{E}\mathbf{A}.$$

Similarly, g_2 can be equivalently written as

$$g_2(\mathbf{X}) = \frac{\tau q}{2} \sum_{\ell=1}^k \|x^\ell\|^2 = \frac{\tau q}{2} \|\mathbf{X}\|_F^2.$$

Hence, g_2 is differentiable and its gradient is given by $\nabla g_2(\mathbf{X}) = \tau q \mathbf{X}$. Therefore,

$$\nabla g(\mathbf{X}) = \nabla g_1(\mathbf{X}) + \nabla g_2(\mathbf{X}) = (m + \tau q)\mathbf{X} - \mathbf{E}\mathbf{A}.$$

Based on the relation (2), finding $\mathbf{X} \in \partial g^*(\mathbf{Y})$ is equivalent to solving the equation

$$\mathbf{Y} = (m + \tau q)\mathbf{X} - \mathbf{E}\mathbf{A}.$$

It follows that

$$\mathbf{X} = \frac{\mathbf{Y} + \mathbf{E}\mathbf{A}}{m + \tau q} \in \partial g^*(\mathbf{Y}).$$

Our goal, then, is to find $\mathbf{Y}_p \in \partial h(\mathbf{X}_p)$ from which we will obtain \mathbf{X}_{p+1} and thereby compute the first N terms of the sequence $\{\mathbf{X}_p\}$ via Algorithm 1. Toward this end we will find subgradients of the convex function h .

For each $i = 1, \dots, m$, let $r(i) \in \{1, \dots, k\}$ be an index for which

$$\sum_{\ell=1, \ell \neq r(i)}^k \|x^\ell - a^i\|^2 = \max_{r=1, \dots, k} \sum_{\ell=1, \ell \neq r}^k \|x^\ell - a^i\|^2,$$

in which case we see that a subgradient $\mathbf{W} \in \partial h_1(\mathbf{X})$ is given by

$$\mathbf{W} = \sum_{i=1}^m \left(\mathbf{X} - \mathbf{A}_i - e_{r(i)}(x^{r(i)} - a^i) \right) = m\mathbf{X} - \mathbf{E}\mathbf{A} - \sum_{i=1}^m e_{r(i)}(x^{r(i)} - a^i), \quad (16)$$

where $\mathbf{A}_i \in \mathbb{R}^{k \times d}$ is the matrix whose all rows are a^i and e_r is the $k \times 1$ column vector with a one in the r^{th} position and zeros elsewhere.

Now for $h_2(\mathbf{X}) = \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \varphi_{\Omega_i^\ell}(x^\ell)$, we have

$$\frac{\partial h_2}{\partial x^j}(\mathbf{X}) = \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \frac{\partial}{\partial x^j} \varphi_{\Omega_i^\ell}(x^\ell) = \tau \sum_{i=1}^q P(x^j; \Omega_i^j)$$

with $j = 1, \dots, k$. Then $\mathbf{U} = 1/\tau \nabla h_2(\mathbf{X})$ is the $k \times d$ matrix whose rows are $u^j = \sum_{i=1}^q P(x^j; \Omega_i^j)$.

The form of the DCA instructs us to find $\mathbf{Y}_p \in \partial h(\mathbf{X}_p)$ at the p th iteration, so we set $\mathbf{Y}_p = \mathbf{W} + \tau \mathbf{U}$. Combining the above results gives $\mathbf{X}_{p+1} = (\mathbf{W} + \tau \mathbf{U} + \mathbf{EA})/(m + \tau q)$. Substituting (16) for \mathbf{W} , we obtain the recursive relation

$$\mathbf{X}_{p+1} = \frac{1}{m + \tau q} \left(m\mathbf{X}_p + \tau \mathbf{U} - \sum_{i=1}^m e_{r(i)} (x_p^{r(i)} - a^i) \right),$$

where x_p^ℓ denotes the ℓ^{th} row of \mathbf{X}_p . The following algorithm summarizes the DCA-based procedure we just derived.

Algorithm 2 DC program for (15)

INPUT: $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau$

for $p = 1, \dots, N$ **do**

for $i = 1, \dots, m$ **do**

 Find $r(i)$ s.t. $\|x_{p-1}^{r(i)} - a^i\|^2 = \min\{\|x_{p-1}^\ell - a^i\|^2 \mid \ell = 1, \dots, k\}$

 Set $\mathbf{W}_i := e_{r(i)} (x_{p-1}^{r(i)} - a^i)$

end

for $\ell = 1, \dots, k$ **do**

 Find $u^\ell := \sum_{j=1}^q P(x_{p-1}^\ell; \Omega_j^\ell)$

end

 Set $\mathbf{X}_p := \frac{1}{m + \tau q} (m\mathbf{X}_{p-1} + \tau \mathbf{U} - \sum_{i=1}^m \mathbf{W}_i)$

end

OUTPUT: \mathbf{X}_N

Inspecting (15), we see that for small τ our problem begins to resemble the associated unconstrained problem. For solving the clustering (14), we may gradually increase the value of the penalty parameter $\tau > 0$ by periodically multiplying by some $\sigma > 1$ and terminate whenever $\tau > \tau_f$. This may be accomplished by Algorithm 3. Notice that for the initial choice of τ , the maximum number of overall iterations of Algorithm 3 is $N \lceil \log_\sigma (\tau_f / \tau) \rceil$, where $\lceil \cdot \rceil$ denotes the ceiling function.

Algorithm 3 Penalty DC program for (14)

INPUT: $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau, \sigma, \tau_f$

while $\tau < \tau_f$ **do**

 Find \mathbf{X}_N by executing **Algorithm 2** with $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau$

 Reassign $\mathbf{X}_0 := \mathbf{X}_N$

 Reassign $\tau := \sigma \tau$

end

OUTPUT: \mathbf{X}_N

Remark 4.1: Assume that in problem (14) the constraint $Q_\ell = \cap_{j=1}^q \Omega_j^\ell$ for each center x^ℓ for $\ell = 1, \dots, k$ is simple enough. Then we can use the *projected k-means algorithm* which is similar to the *k-means algorithm* as follows. After assigning each data point to its nearest centroid, we update each centroid by

$$x^\ell \leftarrow P \left(\frac{\sum_{i \in \mathbf{A}(x^\ell)} a^i}{|\mathbf{A}(x^\ell)|}; Q_\ell \right),$$

where $\mathbf{A}(x^\ell) = \{i \in \{1, \dots, m\} \mid \|a^i - x^\ell\| = \min_{j=1, \dots, k} \|a^i - x^j\|\}$ is the cluster associated with x^ℓ . In this case, to get a good starting point \mathbf{X}_0 for Algorithm 3, we perform the projected *k-means algorithm* in several initial steps. This procedure has shown its efficiency in clustering without constraint; see [1].

5. Set clustering with constraints

In this section, we turn our attention to a model of *set clustering* with constraints, i.e. for given m subsets $\Lambda_1, \dots, \Lambda_m \subset \mathbb{R}^d$, we seek k cluster centers $x^\ell \in \cap_{j=1}^q \Omega_j^\ell$ for $\ell = 1, \dots, k$, where each Ω_j^ℓ is a subset of \mathbb{R}^d . The measurement distance is defined by the squared distance functions to the sets involved. The optimization modeling of the problem to be solved is given by

$$\begin{aligned} \min \quad & \psi(x^1, \dots, x^k) = \sum_{i=1}^m \min_{\ell=1, \dots, k} [d(x^\ell; \Lambda_i)]^2 \\ \text{subject to} \quad & x^\ell \in \cap_{j=1}^q \Omega_j^\ell \quad \text{for } \ell = 1, \dots, k. \end{aligned} \quad (17)$$

Throughout this section, we assume that Λ_i for $i = 1, \dots, m$ and Ω_j^ℓ for $j = 1, \dots, q$ and $\ell = 1, \dots, k$ are nonempty, closed and convex.

Using the penalty method based on distance functions with a parameter $\tau > 0$, we consider the constrained set clustering model:

$$\begin{aligned} \min \quad & f(x^1, \dots, x^k) = \frac{1}{2} \sum_{i=1}^m \min_{\ell=1, \dots, k} [d(x^\ell; \Lambda_i)]^2 + \frac{\tau}{2} \sum_{\ell=1}^k \sum_{j=1}^q [d(x^\ell; \Omega_j^\ell)]^2, \\ & x^1, \dots, x^k \in \mathbb{R}^d. \end{aligned} \quad (18)$$

We will now find a DC decomposition of $f = g - h$ as follows. For each $i = 1, \dots, m$, we have

$$\min_{\ell=1, \dots, k} [d(x^\ell; \Lambda_i)]^2 = \sum_{\ell=1}^k [d(x^\ell; \Lambda_i)]^2 - \max_{r=1, \dots, k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2$$

$$\begin{aligned}
&= \sum_{\ell=1}^k \left(\|x^\ell\|^2 - \varphi_{\Lambda_i}(x^\ell) \right) - \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2 \\
&= \|\mathbf{X}\|_F^2 - \left(\sum_{\ell=1}^k \varphi_{\Lambda_i}(x^\ell) + \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2 \right).
\end{aligned}$$

Furthermore, we have

$$\begin{aligned}
\sum_{\ell=1}^k \sum_{j=1}^q [d(x^\ell; \Omega_j^\ell)]^2 &= \sum_{\ell=1}^k \sum_{j=1}^q \left(\|x^\ell\|^2 - \varphi_{\Omega_j^\ell}(x^\ell) \right) \\
&= q \|\mathbf{X}\|_F^2 - \sum_{\ell=1}^k \sum_{j=1}^q \varphi_{\Omega_j^\ell}(x^\ell).
\end{aligned}$$

Let

$$\begin{aligned}
g_1(\mathbf{X}) &= \frac{m}{2} \|\mathbf{X}\|_F^2, \quad g_2(\mathbf{X}) = \frac{\tau q}{2} \|\mathbf{X}_F\|^2, \\
h_1(\mathbf{X}) &= \sum_{i=1}^m \left(\frac{1}{2} \sum_{\ell=1}^k \varphi_{\Lambda_i}(x^\ell) + \frac{1}{2} \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2 \right), \quad h_2(\mathbf{X}) \\
&= \frac{\tau}{2} \sum_{\ell=1}^k \sum_{j=1}^q \varphi_{\Omega_j^\ell}(x^\ell),
\end{aligned}$$

in which case we have the DC decomposition $f = g - h$, where $g = g_1 + g_2$ and $h = h_1 + h_2$ are convex.

Using the relation (2), we can easily see that $\mathbf{X} = 1/(m + \tau q)\mathbf{Y} \in \partial g^*(\mathbf{Y})$. To apply the DCA from Algorithm 1, we also need to find $\mathbf{Y} \in \partial h(\mathbf{X})$ as $\mathbf{Y} = \mathbf{V} + \mathbf{U}$, where $\mathbf{V} \in \partial h_1(\mathbf{X})$ and $\mathbf{U} \in \partial h_2(\mathbf{X})$.

Now, we focus on finding $\mathbf{V} \in \partial h_1(\mathbf{X})$. Define

$$D_i(\mathbf{X}) = \frac{1}{2} \sum_{\ell=1}^k \varphi_{\Lambda_i}(x^\ell),$$

and

$$F_i(\mathbf{X}) = \frac{1}{2} \max_{r=1,\dots,k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2, \quad i = 1, \dots, m.$$

Then $h_1(\mathbf{X}) = \sum_{i=1}^m [D_i(\mathbf{X}) + F_i(\mathbf{X})]$. Based on Proposition 3.3, we see that $\nabla D_i(\mathbf{X})$ is the $k \times d$ matrix given by

$$\nabla D_i(\mathbf{X}) = \begin{bmatrix} P(x^1; \Lambda_i) \\ \vdots \\ P(x^k; \Lambda_i) \end{bmatrix}.$$

For each $i = 1, \dots, m$, choose an index $r(i)$ such that

$$\max_{r=1, \dots, k} \sum_{\ell=1, \ell \neq r}^k [d(x^\ell; \Lambda_i)]^2 = \sum_{\ell=1, \ell \neq r(i)}^k [d(x^\ell; \Lambda_i)]^2.$$

Now, for $j = 1, \dots, k$, define

$$v_i^j = \begin{cases} x^j - P(x^j; \Lambda_i) & \text{if } j \neq r(i), \\ 0 & \text{if } j = r(i). \end{cases}$$

By (1) and the fact that $\nabla[d(x; \Lambda)]^2 = 2(x - P(x; \Lambda))$ for a nonempty closed convex set Λ , the matrix V_i whose j^{th} row is v_i^j defines a subgradient of F_i at \mathbf{X} . It follows that such a subgradient \mathbf{V} is

$$\mathbf{V} = m\mathbf{X} - \sum_{i=1}^m e_{r(i)} \left(x^{r(i)} - P(x^{r(i)}; \Lambda_i) \right).$$

As computed in the previous section, $\nabla h_2(\mathbf{X}) = \tau \mathbf{U}$, where \mathbf{U} is the $k \times d$ matrix whose ℓ^{th} row is $\sum_{j=1}^k P(x^\ell; \Omega_j^\ell)$ for $\ell = 1, \dots, k$. Consequently, the $k \times d$ matrix

$$\begin{aligned} \mathbf{Y} &= m\mathbf{X} - \sum_{i=1}^m e_{r(i)} \left(x^{r(i)} - P(x^{r(i)}; \Lambda_i) \right) + \tau \mathbf{U} = m\mathbf{X} + \tau \mathbf{U} \\ &\quad - \sum_{i=1}^m e_{r(i)} \left(x^{r(i)} - P(x^{r(i)}; \Lambda_i) \right) \end{aligned}$$

belongs to $\partial h(\mathbf{X})$.

Now, for $p \in N$ such that \mathbf{X}_{p-1} is given, one has

$$\mathbf{Y}_{p-1} = m\mathbf{X}_{p-1} + \tau \mathbf{U}_p - \sum_{i=1}^m e_{r(i)} \left(x_{p-1}^{r(i)} - P(x_{p-1}^{r(i)}; \Lambda_i) \right) \in \partial H(\mathbf{X}_{p-1}),$$

where x_p^ℓ is the ℓ^{th} row of \mathbf{X}_p and \mathbf{U}_p is the $k \times d$ matrix whose ℓ^{th} row is $\sum_{j=1}^k P(x_p^\ell; \Omega_j^\ell)$ for $\ell = 1, \dots, k$. It follows that \mathbf{X}_p from the DCA in Algorithm 1 can be determined by

$$\mathbf{X}_p = \frac{1}{\tau q + m} \left(m\mathbf{X}_{p-1} + \tau \mathbf{U}_p - \sum_{i=1}^m e_{r(i)} \left(x_{p-1}^{r(i)} - P(x_{p-1}^{r(i)}; \Lambda_i) \right) \right).$$

We now adapt Algorithm 4 to solve our set clustering problem. Just as in the previous section, we gradually increase the value of the penalty parameter $\tau > 0$ by periodically multiplying it by some $\sigma > 1$ and stopping when $\tau > \tau_f > 0$. This may be accomplished by Algorithm 5. We again see that for an initial choice of $\tau = \tau_0$, the maximum number of overall iterations of Algorithm 5 is $N \lceil \log_\sigma (\tau_f / \tau_0) \rceil$.

Algorithm 4 DC program for (18)

INPUT: $\mathbf{X}_0, \Lambda_i, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau$
for $p = 1, \dots, N$ **do**
 for $i = 1, \dots, m$ **do**
 for $\ell = 1, \dots, k$ **do**
 Set $w_i^\ell := P(x_{p-1}^\ell; \Lambda_i)$
 end
 Find $r(i)$ s.t. $\|x_{p-1}^{r(i)} - w_i^{r(i)}\|^2 = \min_{\ell=1,\dots,k} \|x_{p-1}^\ell - w_i^\ell\|^2$
 end
 for $\ell = 1, \dots, k$ **do**
 Find $u^\ell := \sum_{j=1}^q P(x_{p-1}^\ell; \Omega_j^\ell)$
 end
 $\mathbf{X}_p := \frac{1}{\tau q + m} (m \mathbf{X}_{p-1} + \tau \mathbf{U}_p - \sum_{i=1}^m e_{r(i)}(x_{p-1}^{r(i)} - w_i^{r(i)}))$
end
OUTPUT: \mathbf{X}_N

Algorithm 5 Penalty DC program for (17)

INPUT: $\mathbf{X}_0, \{\Lambda_i\}_{i=1}^m, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau, \tau_f, \sigma$
while $\tau < \tau_f$ **do**
 Find \mathbf{X}_N by executing **Algorithm 4** with $\mathbf{X}_0, \{\Lambda_i\}_{i=1}^m, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, \tau, N$
 Reassign $\mathbf{X}_0 := \mathbf{X}_N$
 Reassign $\tau := \sigma \tau$
end
OUTPUT: \mathbf{X}_N

6. Multifacility location with constraints

Given a set of m points (nodes) a^1, a^2, \dots, a^m in \mathbb{R}^d , our goal is find k centers x^ℓ for $\ell = 1, \dots, k$, which must be in constraint sets $\bigcap_{i=1}^q \Omega_i^\ell$ for $\ell = 1, \dots, k$, such that the *transportation cost* to the nodes is minimized. The same setting in Section 4 gives us the constrained minimization problem:

$$\min \psi(\mathbf{X}) \quad \text{subject to } \mathbf{X} \in \bigcap_{i=1}^q \Omega_i, \quad (19)$$

where the total cost now is given by

$$\psi(\mathbf{X}) = \psi(x^1, \dots, x^k) = \sum_{i=1}^m \min_{\ell=1,\dots,k} \|x^\ell - a^i\|.$$

This problem can be converted to an unconstrained minimization problem:

$$\begin{aligned} \min f_{\tau}(x^1, \dots, x^k) &= \sum_{i=1}^m \min_{\ell=1, \dots, k} \|x^{\ell} - a^i\| \\ &+ \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q [d(x^{\ell}; \Omega_i^{\ell})]^2, x^1, \dots, x^k \in \mathbb{R}^d. \end{aligned} \quad (20)$$

where $\tau > 0$ is a parameter.

We apply Nesterov's smoothing techniques from [4] to approximate the objective function f_{τ} by a new DC function which is favorable for applying the DCA.

$$\begin{aligned} f_{\tau, \mu}(x^1, \dots, x^k) &= \left(\frac{\mu}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left\| \frac{x^{\ell} - a^i}{\mu} \right\|^2 + \frac{\tau q}{2} \sum_{\ell=1}^k \|x^{\ell}\|^2 \right) \\ &- \left(\frac{\mu}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left[d\left(\frac{x^{\ell} - a^i}{\mu}; \mathbb{B} \right) \right]^2 \right. \\ &\left. + \sum_{i=1}^m \max_{r=1, \dots, k} \sum_{\ell=1, \ell \neq r}^k \|x^{\ell} - a^i\| + \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \varphi_{\Omega_i^{\ell}}(x^{\ell}) \right). \end{aligned}$$

In what follows, we use f instead of $f_{\tau, \mu}$ for the simplicity of notations. The original clustering problem now can be solved using a DC programming:

$$\min f(x^1, \dots, x^k) = g(x^1, \dots, x^k) - h(x^1, \dots, x^k), \quad x^1, \dots, x^k \in \mathbb{R}^d.$$

In this formulation, g and h are convex functions on $(\mathbb{R}^d)^k$ defined by

$$\begin{aligned} g(x^1, \dots, x^k) &= g_1(x^1, \dots, x^k) + g_2(x^1, \dots, x^k), \\ h(x^1, \dots, x^k) &= h_1(x^1, \dots, x^k) + h_2(x^1, \dots, x^k) + h_3(x^1, \dots, x^k), \end{aligned}$$

with their respective components defined as

$$\begin{aligned} g_1 &= \frac{\mu}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left\| \frac{x^{\ell} - a^i}{\mu} \right\|^2, \quad g_2 = \frac{\tau q}{2} \sum_{\ell=1}^k \|x^{\ell}\|^2, \\ h_1 &= \frac{\mu}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left[d\left(\frac{x^{\ell} - a^i}{\mu}; \mathbb{B} \right) \right]^2, \quad h_2 = \sum_{i=1}^m \max_{r=1, \dots, k} \sum_{\ell=1, \ell \neq r}^k \|x^{\ell} - a^i\|, \\ h_3 &= \frac{\tau}{2} \sum_{\ell=1}^k \sum_{i=1}^q \varphi_{\Omega_i^{\ell}}(x^{\ell}). \end{aligned}$$

The function g_1 can be equivalently written as

$$\begin{aligned}
 g_1(\mathbf{X}) &= \frac{1}{2\mu} \sum_{i=1}^m \sum_{\ell=1}^k \|x^\ell - a^i\|^2 \\
 &= \frac{1}{2\mu} \sum_{i=1}^m \sum_{\ell=1}^k \left(\|x^\ell\|^2 - 2\langle x^\ell, a^i \rangle + \|a^i\|^2 \right) \\
 &= \frac{1}{2\mu} \left(m\|\mathbf{X}\|_F^2 - 2\langle \mathbf{X}, \mathbf{EA} \rangle + k\|\mathbf{A}\|_F^2 \right).
 \end{aligned}$$

Note that g_1 is differentiable and its gradient is given by

$$\nabla g_1(\mathbf{X}) = \frac{1}{\mu} [m\mathbf{X} - \mathbf{EA}].$$

The function g_2 is the same as before so its gradient is given by

$$\nabla g_2(\mathbf{X}) = \tau q \mathbf{X}.$$

Since $g(\mathbf{X}) = g_1(\mathbf{X}) + g_2(\mathbf{X})$, its gradient can be computed by

$$\begin{aligned}
 \nabla g(\mathbf{X}) &= \nabla g_1(\mathbf{X}) + \nabla g_2(\mathbf{X}) \\
 &= \frac{1}{\mu} (m\mathbf{X} - \mathbf{EA}) + \tau q \mathbf{X} \\
 &= \left(\frac{m}{\mu} + \tau q \right) \mathbf{X} - \frac{1}{\mu} \mathbf{S},
 \end{aligned}$$

where $\mathbf{S} = \mathbf{EA}$. The latter can equivalently be written as

$$\mathbf{Y} = \left(\frac{m}{\mu} + \tau q \right) \mathbf{X} - \frac{1}{\mu} \mathbf{S}.$$

Our goal now is to compute $\nabla g^*(\mathbf{Y})$, which can be accomplished by the relation (2). Then with some algebraic manipulations, we can show that

$$\nabla g^*(\mathbf{Y}) = \mathbf{X} = \frac{\mu \mathbf{Y} + \mathbf{S}}{m + \mu \tau q}.$$

Next, we will demonstrate in more details the techniques we used in finding a subgradient for the convex function h . Recall that h is defined by

$$h(\mathbf{X}) = \sum_{i=1}^3 h_i(\mathbf{X}).$$

We will start with the function h_1 given by

$$h_1(\mathbf{X}) = \frac{\mu}{2} \sum_{i=1}^m \sum_{\ell=1}^k \left[d\left(\frac{x^\ell - a^i}{\mu}; \mathbb{B}\right) \right]^2.$$

Similar to the situation in [4], we get

$$\frac{\partial h_1}{\partial x^\ell}(\mathbf{X}) = \sum_{i=1}^m \left(\frac{x^\ell - a^i}{\mu} - P\left(\frac{x^\ell - a^i}{\mu}; \mathbb{B}\right) \right).$$

Thus, for $\ell = 1, 2, \dots, k$, $\nabla h_1(\mathbf{X}) = \mathbf{Z}$ is the $k \times d$ matrix whose ℓ^{th} row is $\partial h_1 / \partial x^\ell(\mathbf{X})$.

Let us compute a subgradient of h_2 as in [4]

$$h_2(\mathbf{X}) = \sum_{i=1}^m \max_{\ell=1, \dots, k} \sum_{j=1, j \neq \ell}^k \|x^j - a^i\| = \sum_{i=1}^m \gamma_i(\mathbf{X}),$$

where $\gamma_i(\mathbf{X}) = \max_{\ell=1, \dots, k} \sum_{j=1, j \neq \ell}^k \|x^j - a^i\|$. For each $i = 1, \dots, m$, define

$$\gamma_{i\ell}(\mathbf{X}) = \sum_{j=1, j \neq \ell}^k \|x^j - a^i\|, \quad \ell = 1, \dots, k.$$

Then $\gamma_i(\mathbf{X}) = \max_{\ell=1, \dots, k} \gamma_{i\ell}(\mathbf{X})$.

Based on the subdifferential formula for maximum functions, for each $i = 1, \dots, m$, we find $\mathbf{W}_i \in \partial \gamma_i(\mathbf{X})$. Then define $\mathbf{W} = \sum_{i=1}^m \mathbf{W}_i$ to get a subgradient of the function h_2 at \mathbf{X} by the subdifferential sum rule. To accomplish this goal, we first choose an index $\ell^* = 1, \dots, k$ such that $\gamma_i(\mathbf{X}) = \gamma_{i\ell^*}(\mathbf{X}) = \sum_{j=1, j \neq \ell^*}^k \|x^j - a^i\|$. Using the familiar subdifferential formula of the Euclidean norm function, the j^{th} row w_i^j for $j \neq \ell^*$ of the matrix \mathbf{W}_i is determined as follows

$$w_i^j = \begin{cases} \frac{x^j - a^i}{\|x^j - a^i\|} & \text{if } x^j \neq a^i, \\ 0 & \text{if } x^j = a^i. \end{cases}$$

The ℓ^{*th} row of the matrix \mathbf{W}_i is $w_i^{\ell^*} = 0$.

The procedure for computing $\partial h_3(\mathbf{X})$ is the same in Section 4. Let \mathbf{U} be the matrix whose rows are $\sum_{i=1}^q P(x^\ell; \Omega_i^\ell)$, for $\ell = 1, \dots, k$, then $\nabla h_3(\mathbf{X}) = \tau \mathbf{U}$.

At this point, we are ready to give a new DCA-based algorithm for our problem.

Algorithm 6 DC program for (20)

 INPUT: $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, \tau, \mu, N \in \mathbb{N}$.

for $p = 1, 2, \dots, N$ **do**

 Find $\mathbf{Y}_p := \mathbf{Z}_p + \mathbf{W}_p + \tau \mathbf{U}_p$ where

 $\mathbf{Z}_p := \nabla h_1(\mathbf{X}_p)$

 $\mathbf{W}_p \in \partial h_2(\mathbf{X}_p)$

 $\mathbf{U}_p := \nabla h_3(\mathbf{X}_p)$

 Find $\mathbf{X}_{p+1} := \frac{\mu(\mathbf{Z}_p + \mathbf{W}_p + \tau \mathbf{U}_p) + \mathbf{s}}{m + \mu \tau q}$
end

 OUTPUT: \mathbf{X}_N

We also present below an adapted version of Algorithm 6 for solving (19). We may improve Algorithm 6 by gradually increasing and decreasing the value of the penalty parameter τ and the smoothing parameter μ respectively. This can be done by periodically multiplying them by some $\sigma > 1$, $0 < \delta < 1$ and stopping when $\tau > \tau_f$, $\mu < \mu_f$.

Algorithm 7 Penalty DC program for (19)

 INPUT: $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau, \sigma, \tau_f, \mu, \delta, \mu_f$
while $\tau < \tau_f$ and $\mu > \mu_f$ **do**

 Find \mathbf{X}_N by executing **Algorithm 6** with $\mathbf{A}, \mathbf{X}_0, \{\Omega_j^\ell\}_{j=1,\dots,q}^{\ell=1,\dots,k}, N, \tau, \mu$

 Reassign $\mathbf{X}_0 := \mathbf{X}_N$

 Reassign $\tau := \sigma \tau$

 Reassign $\mu := \delta \mu$
end

 OUTPUT: \mathbf{X}_N

7. Numerical experiments

Example 7.1: We now consider the dataset EIL76 taken from the Traveling Salesman Problem Library [13]. We impose the following constraints on the solution:

- (1) The first center is a common point of a box whose vertices are (40, 40); (40, 60); (20, 60); (20, 40) and a ball of radius $r = 7$ centered at (20, 60).
- (2) The second center is in the intersection of two balls of the same radius $r = 7$, centered at (35, 20) and (45, 22), respectively.

Choosing $\tau = 1$, $\sigma = 10$, $\tau_f = 10^8$, Algorithm 3 yields an approximate solution:

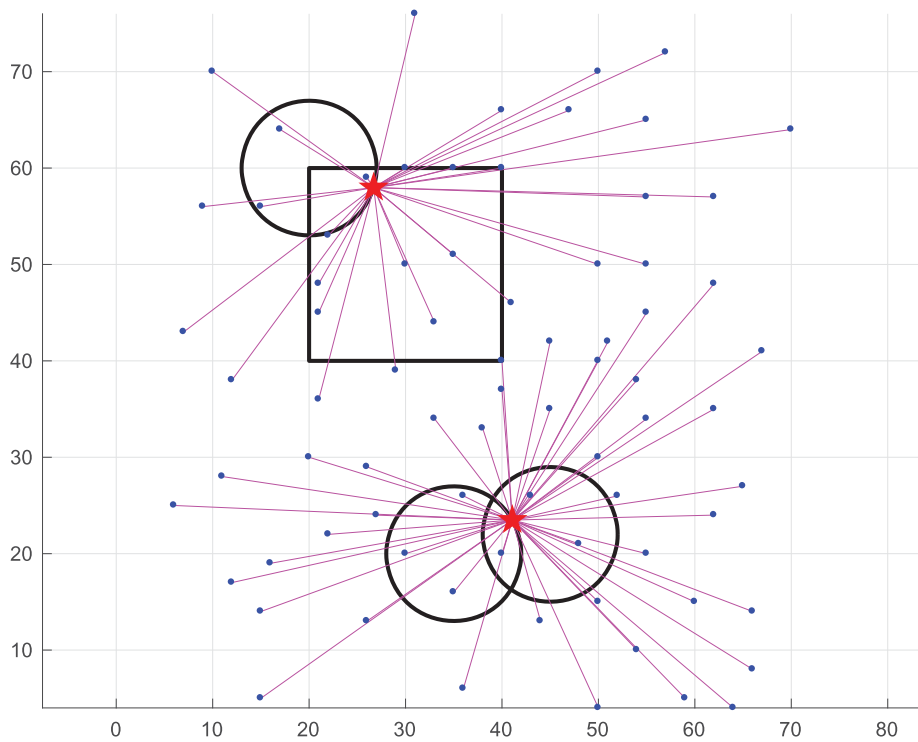


Figure 1. A 2-center constrained clustering problem for dataset EIL76.

$$\mathbf{X} = \begin{pmatrix} 26.69959 & 57.97125 \\ 41.06910 & 23.48799 \end{pmatrix}, \text{ with the cost } \psi(\mathbf{X}) = 33576.25387; \text{ see Figure 1.}$$

7.1. Set clustering with constraints

Example 7.2: We now use Algorithm 5 to solve a set clustering problem with constraints. We consider the latitude and longitude of the 50 most populous US cities taken from 2014 United States Census Bureau data¹, and approximate each city by a ball with radius $0.1\sqrt{A/\pi}$ where A is the city's reported area in square miles.

We use Algorithm 5 for solving 3-center problem generated by this 50-set dataset with requirement that each center must belong to the intersection of two balls. The centers of these constrained balls are the columns of the matrix below

$$\begin{pmatrix} -80 & -80 & -92 & -90 & -115 & -110 \\ 34 & 38 & 37 & 40 & 45 & 40 \end{pmatrix}$$

with corresponding radii given by $(2 \ 3 \ 4 \ 3 \ 4 \ 4)$. The result is plotted in Figure 2 using a plate Carrée projection².

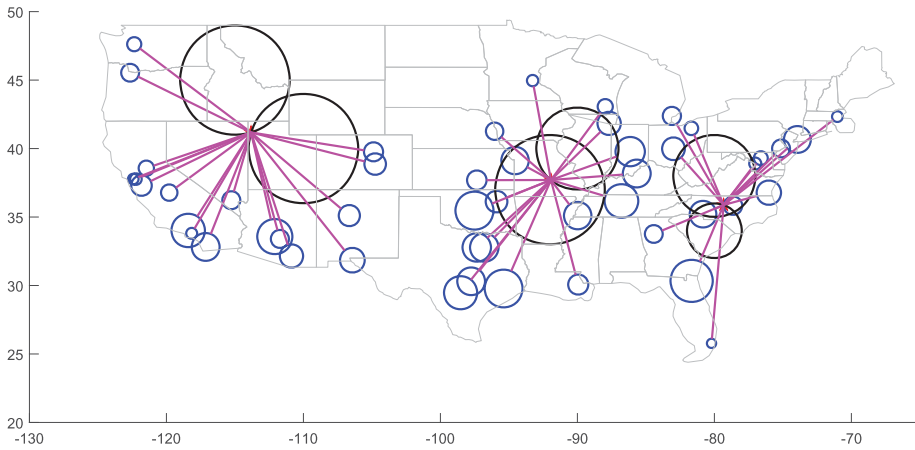


Figure 2. A 3-center set clustering problems with 50 most populous US cities. Each city is approximated by a ball proportional to its area.

We again choose $\tau = 1$, $\sigma = 10$, $\tau_f = 10^8$, Algorithm 5 yields an approximate optimal value $\psi(\mathbf{X}) = 2271.09657$ at an approximate solution given by

$$\mathbf{X} = \begin{pmatrix} -79.32172 & 35.88148 \\ -91.93134 & 37.70436 \\ -113.82289 & 41.17711 \end{pmatrix}.$$

7.2. Multifacility location with constraints

Example 7.3: We now test Algorithm 7 on a data set \mathbf{A} containing random points in 4 balls of radius $r = 0.3$ centered at $(2, 2)$, $(4, 2)$, $(4, 4)$ and $(2, 4)$. Let $k = 4$ and the constraint be the ball with the same radius, centered at $(3, 3)$. We use the *kmeans* (a MATLAB built in function) to partition the nodes into 4 clusters first, and then we selected the 4 cluster centroid locations as starting centers. We choose $\tau = 1$, $\sigma = 10$, $\tau_f = 10^8$, $\mu = 1$, $\delta = 0.75$, $\mu_f = 10^{-6}$. Typical centers are the intersections of the constraint ball boundary and the line connecting centers of each ball to the center of the constraint one. A visualization is shown in Figure 3.

Example 7.4: Next we consider the latitude and longitude data of the $m = 988$ most-populated cities in the contiguous 48 United States [14]. We impose the following constraints on the solution:

- (1) One center is to lie west of -115° longitude and within 4° latitude/longitude of Caldwell, Idaho.
- (2) One center is to lie within the state of Colorado and within 6° latitude/longitude of Oklahoma City, Oklahoma.

- (3) One center is to lie within 2° latitude/longitude of Skokie, Illinois and the triangle with vertices at Cleveland, Ohio; Atlanta, Georgia; and Des Moines, Iowa.
- (4) One center is to lie within 4° latitude/longitude of New York, NY and Washington, DC.

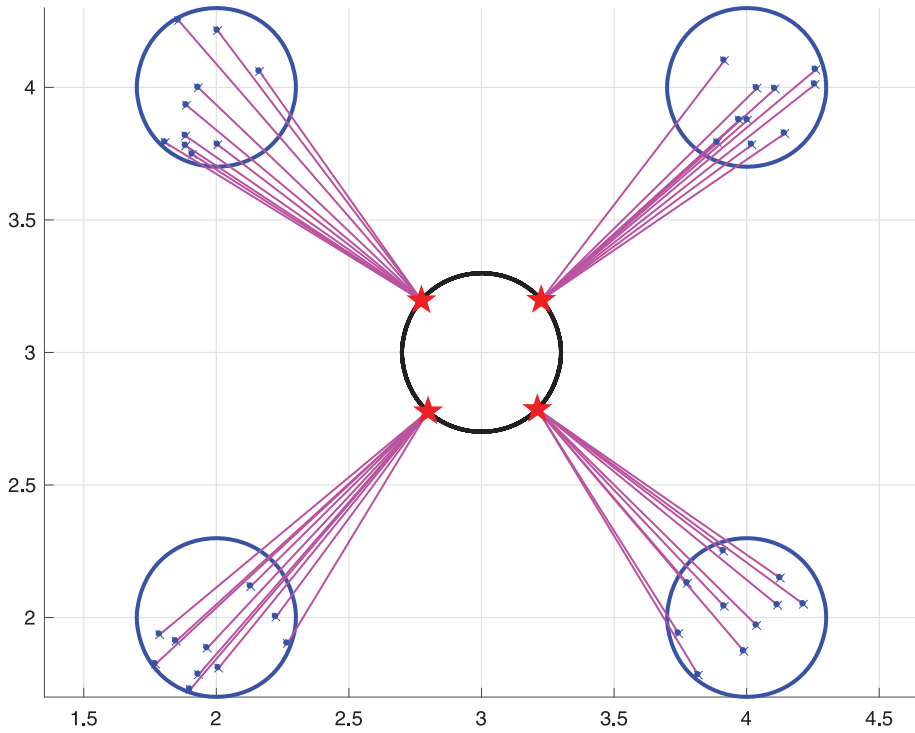


Figure 3. A 4-center multifacility location with one ball constraint.

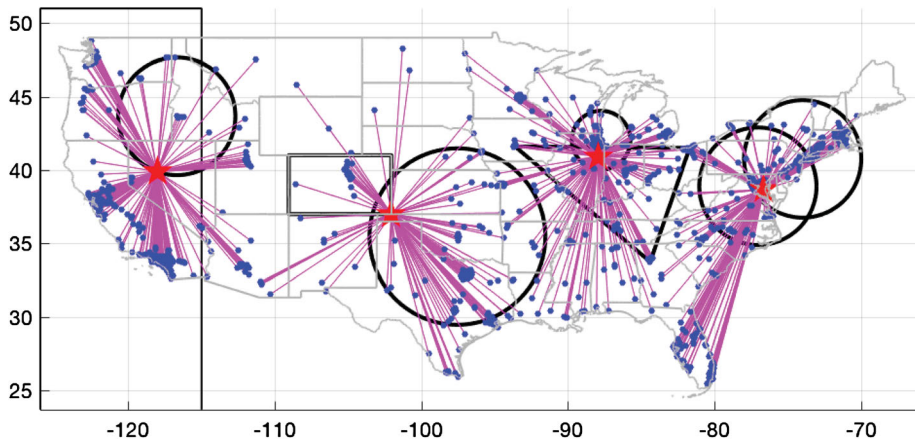


Figure 4. A 4-center constrained multifacility location problem with US cities dataset.

Employing Algorithms 6 and 7 with $\tau = 1$, $\sigma = 100$, $\tau_f = 10^8$, $\mu = 1$, $\delta = 0.85$, $\mu_f = 10^{-6}$, we terminate when $\|\mathbf{X}_{p+1} - \mathbf{X}_p\|_F < 10^{-6}$ and find final centers at

$$\mathbf{X} = \begin{pmatrix} -118.03185 & 39.89550 \\ -102.04996 & 36.99996 \\ -87.93854 & 40.90443 \\ -76.63980 & 38.67968 \end{pmatrix}$$

with an objective value $\psi(\mathbf{X}) = 42586.65060$; see Figure 4.

Notes

1. https://en.wikipedia.org/wiki/List_of_United_States_cities_by_population
2. <https://www.mathworks.com/help/map/pcarree.html>

Disclosure statement

No potential conflict of interest was reported by the authors.

Funding

Research of Nguyen Mau Nam (Mau Nguyen) was partly supported by the National Science Foundation under grant DMS-1716057. Research of Nguyen Thai An was supported by the Vietnam National Foundation for Science and Technology Development under grant #101.01-2017.325

References

- [1] An LTH, Belghiti MT, Tao PD. A new efficient algorithm based on DC programming and DCA for clustering. *J Glob Optim.* 2007;27:503–608.
- [2] An LTH, Minh LH, Tao PD. New and efficient DCA based algorithms for minimum sum-of-squares clustering. *Pattern Recognit.* 2014;47:388–401.
- [3] Nam NM, Geremew W, Reynolds S, et al. The nesterov smoothing technique and minimizing differences of convex functions for hierarchical clustering. *Optim Lett.* 2018;12:455–473. submitted.
- [4] Nam NM, Rector RB, Giles D. Minimizing differences of convex functions with applications to facility location and clustering. *J Optim Theory Appl.* 2017;173:255–278. submitted.
- [5] Hiriart-Urruty JB, Lemaréchal C. *Fundamental of convex analysis*. Berlin: Springer-Verlag; 2001.
- [6] Mordukhovich BS, Nam NM. *An easy path to convex analysis and applications*. San Rafael, CA: Morgan & Claypool Publishers; 2014.
- [7] Rockafellar RT. *Convex analysis*. Princeton, NJ: Princeton University Press; 1970.
- [8] Tao PD, An LTH. Convex analysis approach to D.C. programming: theory algorithms and applications. *Acta Math Vietnam.* 1997;22:289–355.
- [9] Tao PD, An LTH. A D.C. optimization algorithm for solving the trust-region subproblem. *SIAM J Optim.* 1998;8:476–505.
- [10] Chi E, Zhou H, Lange K. Distance majorization and its applications. *Math Program Series A.* 2014;146:409–436.

- [11] Nocedal J, Wright SJ. Numerical optimization. 2nd Edition. New York: Springer; 2006.
- [12] Mordukhovich BS, Nam NM. Geometric approach to convex subdifferential calculus. *Optim.* 2017;66:839–873.
- [13] Reinelt G. TSPLIB: A traveling salesman problem library. *ORSA J Comput.* 1991;3: 376–384.
- [14] United States Cities Database. Simple Maps: Geographic Data Products, 2017, <http://simplemaps.com/data/us-cities>.