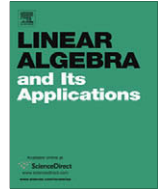




Contents lists available at ScienceDirect

Linear Algebra and its Applications

journal homepage: www.elsevier.com/locate/laa



Contour approximation of data: A duality theory

Cem Iyigun*, Adi Ben-Israel

RUTCOR – Rutgers Center for Operations Research, Rutgers University, 640 Bartholomew Rd., Piscataway, NJ 08854-8003, USA

ARTICLE INFO

Article history:

Received 6 December 2007

Accepted 6 January 2009

Submitted by R.A. Brualdi

Keywords:

Clustering

Contour approximation of data

Duality

Mahalanobis distance

Harmonic mean

Joint distance function

Weiszfeld method

ABSTRACT

Given a dataset \mathcal{D} partitioned in clusters, the joint distance function (JDF) $J(\mathbf{x})$ at any point \mathbf{x} is the harmonic mean of the distances between \mathbf{x} and the cluster centers. The JDF is a continuous function, capturing the data points in its lower level sets (a property called contour approximation), and is a useful concept in probabilistic clustering and data analysis.

In particular, contour approximation allows a compact representation of the data: for a dataset in \mathbb{R}^n with N points organized in K clusters, the JDF requires K centers and covariances (if Mahalanobis distances are used), for a total of $Kn(n+3)/2$ parameters, and a considerable reduction of storage if $N \gg K, n$.

The JDF of the whole dataset, $J(\mathcal{D}) := \sum \{J(\mathbf{x}) : \mathbf{x} \in \mathcal{D}\}$, is a measure of the classifiability of the data, and can be used to determine the “right” number of clusters for \mathcal{D} .

A duality theory for the JDF $J(\mathcal{D})$ is given, in analogy with Kuhn’s geometric duality theory for the Fermat–Weber location problem. The JDF $J(\mathcal{D})$ is the optimal value of a primal problem (P), for which a dual problem (D) is given, with a sharp lower bound on $J(\mathcal{D})$.

© 2009 Elsevier Inc. All rights reserved.

1. Introduction

We use the abbreviation

$$\overline{1, K} := \{1, 2, \dots, K\} \tag{1}$$

for the indicated index set. The standard inner product in \mathbb{R}^n is denoted by $\mathbf{x} \cdot \mathbf{y}$, and for a positive definite matrix Q , the **elliptic norm**,

* Corresponding author.

E-mail addresses: iyigun@rutcor.rutgers.edu (C. Iyigun), adi.benisrael@gmail.com (A. Ben-Israel).

$$\|\mathbf{u}\|_Q := (\mathbf{u} \cdot \mathbf{Q} \mathbf{u})^{1/2}. \tag{2}$$

The **Euclidean norm**

$$\|\mathbf{u}\| := (\mathbf{u} \cdot \mathbf{u})^{1/2} \tag{3}$$

corresponds to $Q = I$, in which case the subscript is omitted. We note the relation

$$\|\mathbf{u}\|_Q = \|Q^{1/2} \mathbf{u}\|, \quad \forall \mathbf{u} \in \mathbb{R}^n. \tag{4}$$

We take data points $\mathbf{x} = (x_1, \dots, x_n)$ as vectors in \mathbb{R}^n . Let $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^n$ be a dataset with N points, partitioned into K **clusters**

$$\mathcal{D} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \dots \cup \mathcal{C}_K, \quad \text{where } \mathcal{C}_i \cap \mathcal{C}_j = \emptyset, \quad \text{if } i \neq j.$$

The k th cluster \mathcal{C}_k has a **center** \mathbf{c}_k and an associated **distance function** $d_k(\mathbf{x}, \mathbf{c}_k)$, defined by

$$d_k(\mathbf{x}, \mathbf{c}_k) := \|\mathbf{x} - \mathbf{c}_k\|_{Q_k}, \text{ occasionally abbreviated by } d_k(\mathbf{x}), \tag{5}$$

where the positive definite matrix Q_k models the geometry of the cluster. In particular, the **Mahalanobis distance**

$$d_k(\mathbf{x}, \mathbf{c}_k) := \sqrt{(\mathbf{x} - \mathbf{c}_k) \cdot \Sigma_k^{-1} (\mathbf{x} - \mathbf{c}_k)}, \tag{6}$$

where Σ_k is the covariance matrix of the data involved.

In probabilistic clustering, for each point \mathbf{x} there are probabilities,

$$p_k(\mathbf{x}) = \text{Prob}\{\mathbf{x} \in \mathcal{C}_k\}, \quad k \in \overline{1, K}, \tag{7}$$

of belonging to the K clusters. These membership probabilities are assumed to depend on the distances $\{d_k(\mathbf{x})\}$ between \mathbf{x} and the cluster centers, with cluster membership more probable the closer is the cluster center. A simple such model is

$$p_k(\mathbf{x}) d_k(\mathbf{x}) = J(\mathbf{x}), \quad \forall k \in \overline{1, K}, \tag{8}$$

where $J(\mathbf{x})$ is a function of \mathbf{x} that does not depend on k , see [3]. In what follows we use (8) as our working principle.

Since probabilities add to one, (8) gives

$$J(\mathbf{x}) = \frac{\prod_{k=1}^K d_k(\mathbf{x})}{\sum_{k=1}^K \prod_{j \neq k} d_j(\mathbf{x})}, \tag{9}$$

which is (up to the constant K) the harmonic mean of the K distances $d_k(\mathbf{x})$. We call $J(\cdot)$ the **joint distance function** (or JDF for short). $J(\mathbf{x})$ has the dimension of distance, and is an indicator of the **classifiability** of the point \mathbf{x} , which is easier to classify the smaller is $J(\mathbf{x})$. In particular, $J(\mathbf{x}) = 0$ if and only if \mathbf{x} coincides with one of the centers \mathbf{c}_k , in which case $p_k(\mathbf{x}) = 1$ (and $p_i(\mathbf{x}) = 0$ for $i \neq k$).

Since $J(\mathbf{x}) = (\sum_k p_k(\mathbf{x})) J(\mathbf{x}) = \sum_k p_k(\mathbf{x}) (p_k(\mathbf{x}) d_k(\mathbf{x}))$ it follows that:

$$J(\mathbf{x}) = \sum_{k=1}^K p_k(\mathbf{x})^2 d_k(\mathbf{x}), \tag{10}$$

an alternative expression of the JDF.

The JDF has an important approximation property: it captures the data points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ in its lower level sets. This property, called **contour approximation**, was studied in [2] where the significance of the harmonic mean was elucidated. See Fig. 1 for an illustration of contour approximation for datasets with 3 clusters in \mathbb{R}^2 .

The JDF of the whole dataset \mathcal{D} is defined as the sum over all data points,

$$\begin{aligned} J(\mathcal{D}) &= \sum_{i=1}^N J(\mathbf{x}_i) \\ &= \sum_{k=1}^K \sum_{i=1}^N p_k(\mathbf{x}_i)^2 d_k(\mathbf{x}_i), \text{ by (10)}, \end{aligned} \tag{11}$$

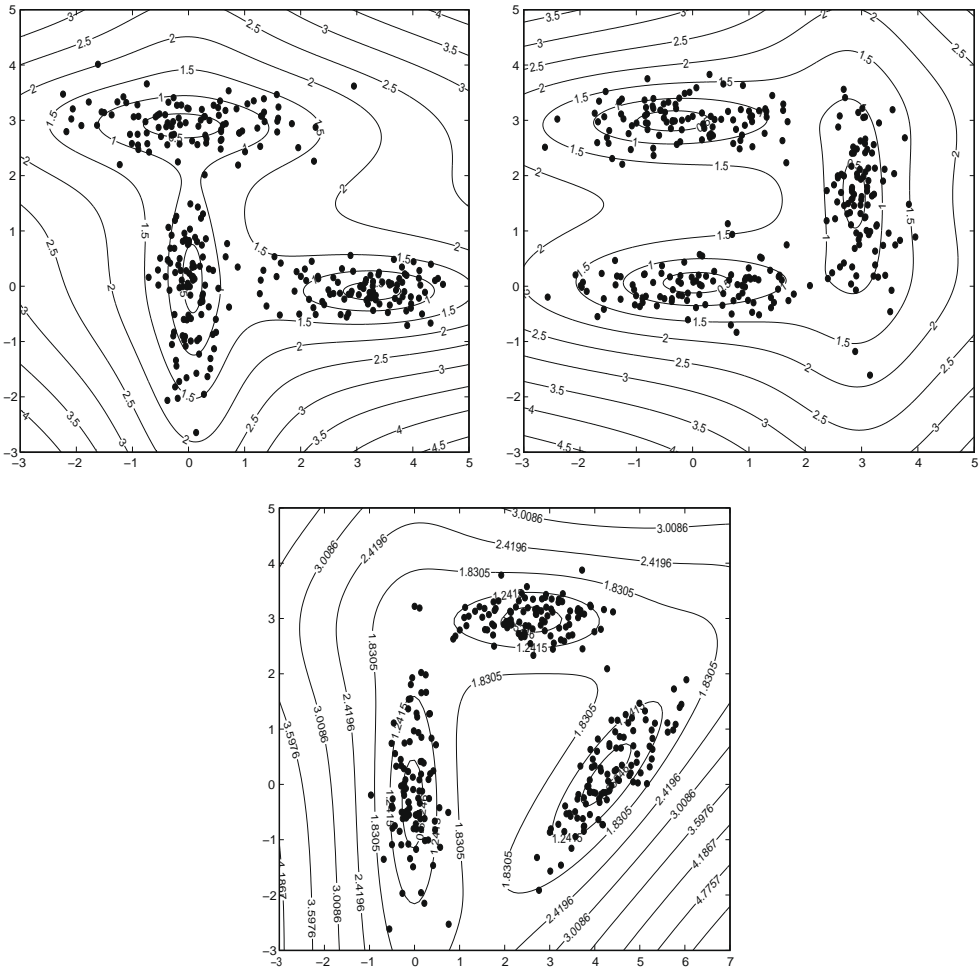


Fig. 1. Contour approximation of data by the joint distance function.

measuring the uncertainty of classifying the dataset \mathcal{D} . This suggests the following formulation of clustering as a minimization problem,

$$\begin{aligned}
 \min \quad & \sum_{k=1}^K \sum_{i=1}^N p_k(\mathbf{x}_i)^2 d_k(\mathbf{x}_i, \mathbf{c}_k) & (P) \\
 \text{s.t.} \quad & \sum_{k=1}^K p_k(\mathbf{x}_i) = 1, \quad i \in \overline{1, N}, \\
 & p_k(\mathbf{x}_i) \geq 0, \quad k \in \overline{1, K}, \quad i \in \overline{1, N},
 \end{aligned}$$

called the **primal problem**, with two sets of variables, the **centers** $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ and the **probabilities** $\{p_k(\mathbf{x}_i) : k \in \overline{1, K}, i \in \overline{1, N}\}$. A natural approach is to fix one set of variables, and minimize (P) with respect to the other set, then fix the other set, etc. We thus alternate between

- (1) the **probabilities problem**, that is (P) with given centers, and
- (2) the **centers problem**, (P) with given probabilities.

This is the essence of the **probabilistic distance clustering** method of [3].

Plan: We present a duality theory for contour approximation, that is based on the geometric duality developed by Kuhn [7] for the Fermat–Weber location problem.

In Section 2, we outline the algorithm of [3] for an iterative solution of (P). The problem encountered when one of the centers coincides with a data point is addressed in Section 3, where a modified gradient is constructed, and applied in Theorem 1 to characterize optimality in the centers problem. In Section 4, we introduce a dual problem (D) for the centers problem, and prove weak duality, Theorem 2. Strong duality is established in Section 5, Theorems 3 and 4, in that the dual pair $\{(P),(D)\}$ have no duality gap.

Notes

- (a) Taking data points as vectors in \mathbb{R}^n assumes continuous data. Indeed, discrete data is not preserved by vector operations (for example, the sum $\mathbf{x} + \mathbf{y}$ of two data points is not necessarily a data point).
- (b) The presence of squares of probabilities in (P) is explained as a device for smoothing the underlying clustering problem which is non-smooth. For related results in greater generality see the seminal paper [11].
- (c) Eq. (8) are the optimality conditions for the probabilities problem (P).
- (d) Contour approximation allows a compact representation of the data in question. Consider for example a dataset \mathcal{D} in \mathbb{R}^n with N points that is organized in K clusters. Then the JDF $J(\mathcal{D})$ of (9) requires K centers and K covariances (if Mahalanobis distances are used), for a total of $Kn(n+3)/2$ parameters, a considerable saving if $N \gg K, n$.

2. Probabilistic distance clustering

Recall that problem (P) has two sets of variables, centers $\{\mathbf{c}_k\}$ and probabilities $\{p_k(\mathbf{x}_i)\}$. We present updates for the probabilities problem and the centers problem.

Probabilities update: The centers $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ are assumed given, and the distances $d_k(\mathbf{x}_i)$ are computed for all centers \mathbf{c}_k and data points \mathbf{x}_i . The minimizing probabilities are explicitly computed from equations (8) as follows:

$$p_k(\mathbf{x}_i) = \frac{\prod_{j \neq k} d_j(\mathbf{x}_i)}{\sum_{m=1}^K \prod_{j \neq m} d_j(\mathbf{x}_i)}, \quad k \in \overline{1, K}. \tag{12}$$

Centers update: Fixing the probabilities $p_k(\mathbf{x}_i)$ in (P), the objective function is a function of the cluster centers,

$$f(\mathbf{c}_1, \dots, \mathbf{c}_K) = \sum_{k=1}^K \sum_{i=1}^N p_k(\mathbf{x}_i)^2 d_k(\mathbf{x}_i, \mathbf{c}_k), \tag{13}$$

where $d_k(\mathbf{x}_i, \mathbf{c}_k) = \|\mathbf{x}_i - \mathbf{c}_k\|_{Q_k}$, an elliptic distance defined by a positive definite matrix Q_k that is associated with the k th cluster. The gradient of (13) with respect to \mathbf{c}_k , at a variable point \mathbf{c} , is

$$\nabla_{\mathbf{c}_k} f(\mathbf{c}) = -Q_k \sum_{i=1}^N \frac{p_k(\mathbf{x}_i)^2}{d_k(\mathbf{x}_i, \mathbf{c})} (\mathbf{x}_i - \mathbf{c}). \tag{14}$$

Zeroing the gradients (14) we get the optimal centers $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ as convex combinations of the data points,

$$\mathbf{c}_k = \sum_{i=1}^N \lambda_k(\mathbf{x}_i) \mathbf{x}_i, \tag{15a}$$

where the weights $\lambda_k(\mathbf{x}_i)$ are given by

$$\lambda_k(\mathbf{x}_i) = \frac{p_k(\mathbf{x}_i)^2 / d_k(\mathbf{x}_i, \mathbf{c}_k)}{\sum_{j=1}^N p_k(\mathbf{x}_j)^2 / d_k(\mathbf{x}_j, \mathbf{c}_k)}, \quad k \in \overline{1, K}, \quad i \in \overline{1, N}. \tag{15b}$$

Notes

- (a) The results (15) are circular, in that the centers \mathbf{c}_k are given by weights that depend on the centers. Still, (15a) and (15b) are useful as an iterative method for updating the centers, which is a generalization to several centers of the **Weiszfeld method** [13] for solving the Fermat–Weber location problem.
- (b) Substituting (12) in (15b) shows that the centers update can be done in terms of the distances $\{d_k(\mathbf{x}_i)\}$ alone, and that the probabilities $\{p_k(\mathbf{x}_i)\}$ are not explicitly needed in the computations.
- (c) If Mahalanobis distances (6) are used, it is required to update the estimates of the covariance matrices $\{\Sigma_k\}$, in addition to the updates of probabilities and centers.
- (d) There are cases where the cluster sizes are unknowns to be estimated. For example, data sampled from a distribution that is itself a mixture of several distributions, where it is required to estimate the weights of the mixture, as well as the parameters of the distributions in the mix. The probabilistic distance clustering method of [3] was adapted in [5] to handle such applications.

3. The modified gradient

The gradient (14) is undefined (0/0) if \mathbf{c} coincides with any of the data points \mathbf{x}_i . If a center \mathbf{c}_k coincides with a data point \mathbf{x}_j then \mathbf{x}_j belongs with certainty to the k th cluster and, by (12),

$$p_k(\mathbf{x}_j) = 1, \quad p_m(\mathbf{x}_j) = 0 \quad \text{for all } m \neq k. \tag{16}$$

In this case, we modify the gradient, following [7,8], and denote the modified gradient by $-\mathbf{R}_k$.

If a center \mathbf{c}_k is not a data point, copy (14) with a change of sign,

$$\mathbf{R}_k(\mathbf{c}_k) := Q_k \sum_{i=1}^N \frac{p_k(\mathbf{x}_i)^2}{d_k(\mathbf{x}_i, \mathbf{c}_k)} (\mathbf{x}_i - \mathbf{c}_k). \tag{17a}$$

Otherwise, if a center \mathbf{c}_k coincides with a data point \mathbf{x}_j , define,

$$\mathbf{R}_k(\mathbf{x}_j) := \max \left\{ \|Q_k^{-1/2} \mathbf{R}_k^j\| - p_k(\mathbf{x}_j)^2, 0 \right\} \frac{\mathbf{R}_k^j}{\|\mathbf{R}_k^j\|}, \tag{17b}$$

where

$$\mathbf{R}_k^j = Q_k \sum_{i \neq j} \frac{p_k(\mathbf{x}_i)^2}{d_k(\mathbf{x}_i, \mathbf{x}_j)} (\mathbf{x}_i - \mathbf{x}_j), \tag{17c}$$

and by (16), $p_k(\mathbf{x}_j) = 1$ here and throughout this section.

In (17b), if $\|Q_k^{-1/2} \mathbf{R}_k^j\| < 1$ then $\mathbf{R}_k(\mathbf{x}_j) = \mathbf{0}$; otherwise, $\mathbf{R}_k(\mathbf{x}_j)$ is a vector with magnitude $\|Q_k^{-1/2} \mathbf{R}_k^j\| - 1$ and direction \mathbf{R}_k^j .

Next, a characterization of optimality in terms of the modified gradient.

Theorem 1. *Given the data $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, let $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ be any K points, and let the corresponding probabilities $\{p_k(\mathbf{x}_i) : k \in \overline{1, K}, i \in \overline{1, N}\}$ be given by (12). Then the condition*

$$\mathbf{R}_k(\mathbf{c}_k) = \mathbf{0}, \quad \text{for all } k \in \overline{1, K} \tag{18}$$

is necessary and sufficient for the points $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ to minimize the function f of (13).

Proof. If \mathbf{c}_k is not one of the data points, then $-\mathbf{R}_k(\mathbf{c}_k)$ is the gradient (14) at \mathbf{c}_k , and (18) is both necessary and sufficient for a minimum, by the convexity of f (as a function of \mathbf{c}_k).

If \mathbf{c}_k coincides with a data point \mathbf{x}_j , consider the change from \mathbf{x}_j to $\mathbf{x}_j + t \mathbf{z}$ where $\|\mathbf{z}\|_{Q_k} = 1$. Then,

$$\left. \frac{d}{dt} f(\mathbf{c}_1, \dots, \mathbf{c}_{k-1}, \mathbf{x}_j + t \mathbf{z}, \mathbf{c}_{k+1}, \dots, \mathbf{c}_K) \right|_{t=0} = p_k(\mathbf{x}_j)^2 - \mathbf{R}_k^j \cdot \mathbf{z}. \tag{19}$$

The greatest decrease of (19) is for \mathbf{z} along \mathbf{R}_k^j , i.e., when

$$\mathbf{z} = \frac{\mathbf{R}_k^j}{\|\mathbf{R}_k^j\|_{Q_k}}.$$

Therefore, \mathbf{c}_k (that coincides with \mathbf{x}_j) is a local minimum if and only if,

$$p_k(\mathbf{x}_j)^2 - \frac{\mathbf{R}_k^j \cdot \mathbf{R}_k^j}{\|\mathbf{R}_k^j\|_{Q_k}} \geq 0 \quad \text{or} \quad \frac{(Q_k^{1/2} \mathbf{R}_k^j) \cdot (Q_k^{-1/2} \mathbf{R}_k^j)}{\|Q_k^{1/2} \mathbf{R}_k^j\|} \leq p_k(\mathbf{x}_j)^2,$$

which is equivalent to

$$\|Q_k^{-1/2} \mathbf{R}_k^j\| \leq p_k(\mathbf{x}_j)^2$$

or $\mathbf{R}_k(\mathbf{c}_k) = \mathbf{0}$, by (17b). \square

4. The dual problem

We abbreviate the probabilities $p_k(\mathbf{x}_i)$ by p_{ki} , for $k \in \overline{1, K}$, $i \in \overline{1, N}$. A dual problem (D) for (P) is now given. It uses the data

$$\mathcal{S} := \{\mathbf{x}_i : i \in \overline{1, N}\}, \quad \{p_{ki} : k \in \overline{1, K}, i \in \overline{1, N}\}, \quad \{Q_k : k \in \overline{1, K}\}, \tag{20}$$

consisting of the data points $\{\mathbf{x}_i\}$, the probabilities $\{p_{ki}\}$, and the matrices $\{Q_k\}$ used in the elliptic distances. The dual variables are KN vectors $\{\mathbf{u}_{ki} : k \in \overline{1, K}, i \in \overline{1, N}\}$, one for each cluster and data point. We denote the set of dual variables by \mathbf{U} .

The **dual problem** is:

$$\max \quad g(\mathbf{U}) = \sum_{k=1}^K \sum_{i=1}^N \mathbf{u}_{ki} \cdot \mathbf{x}_i \tag{D}$$

$$\text{s.t.} \quad \sum_{i=1}^N \mathbf{u}_{ki} = \mathbf{0}, \quad k \in \overline{1, K}, \tag{21}$$

$$\|Q_k^{-1/2} \mathbf{u}_{ki}\| \leq p_{ki}^2, \quad i \in \overline{1, N}, \quad k \in \overline{1, K}. \tag{22}$$

Variables $\mathbf{U} = \{\mathbf{u}_{ki}\}$ satisfying (21) and (22) are called **feasible**.

Problem (D) is a generalization of the dual problem for the single facility location problem, see [7,4, Section 1.1.2].

Theorem 2. Let the data \mathcal{S} in (20) be given. Then for any set of centers $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$, and any set of feasible dual variables $\mathbf{U} = \{\mathbf{u}_{ki}\}$,

$$g(\mathbf{U}) \leq f(\mathbf{c}_1, \dots, \mathbf{c}_K). \tag{23}$$

Proof. The objective function $g(\mathbf{U})$ can be written, using (21), as

$$g(\mathbf{U}) = \sum_{k=1}^K \left(\sum_{i=1}^N \mathbf{u}_{ki} \cdot \mathbf{x}_i - \left(\sum_{i=1}^N \mathbf{u}_{ki} \right) \cdot \mathbf{c}_k \right) = \sum_{k=1}^K \sum_{i=1}^N \mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k). \tag{24}$$

Using

$$\mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k) = (Q_k^{-1/2} \mathbf{u}_{ki}) \cdot (Q_k^{1/2} (\mathbf{x}_i - \mathbf{c}_k)),$$

we get from the Cauchy–Schwartz inequality,

$$\begin{aligned}
 |\mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k)| &= \left| Q_k^{-1/2} \mathbf{u}_{ki} \cdot Q_k^{1/2} (\mathbf{x}_i - \mathbf{c}_k) \right| \\
 &\leq \|Q_k^{-1/2} \mathbf{u}_{ki}\| \|\mathbf{x}_i - \mathbf{c}_k\|_{Q_k}, \\
 \therefore g(\mathbf{U}) &= \sum_{k=1}^K \sum_{i=1}^N \mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k) \leq \sum_{k=1}^K \sum_{i=1}^N \|Q_k^{-1/2} \mathbf{u}_{ki}\| \|\mathbf{x}_i - \mathbf{c}_k\|_{Q_k}
 \end{aligned} \tag{25a}$$

$$\begin{aligned}
 &\leq \sum_{k=1}^K \sum_{i=1}^N p_{ki}^2 d_k(\mathbf{x}_i, \mathbf{c}_k), \text{ by (22)} \\
 &= f(\mathbf{c}_1, \dots, \mathbf{c}_K). \quad \square
 \end{aligned} \tag{25b}$$

Note. The dual problem (D) admits a simple mechanical model that we now describe. Consider the space as a “horizontal” board with holes drilled in the N locations $\{\mathbf{x}_i\}$. For each $k \in \overline{1, K}$, let N (weightless, frictionless, zero thickness) strings be tied together in one knot, and let the loose ends pass through the N holes, each attached to a weight. The string from the k th knot through \mathbf{x}_i is connected to a weight of magnitude p_{ki}^2 . The K knots are free to move, and their locations are denoted by $\{\mathbf{c}_k\}$.

The dual variables are interpreted as forces due to the the weights, with \mathbf{u}_{ki} the (negative of the) force exerted on the k th knot \mathbf{c}_k by the i th weight p_{ki}^2 . The K knots will come to stop, because the resultant forces are zero, by (21). For the “right” probabilities, the knots will stop at the optimal centers.

The optimal forces $\{\mathbf{u}_{ki}\}$ minimize the negative of (24),

$$\sum_{k=1}^K \sum_{i=1}^N \mathbf{u}_{ki} \cdot (\mathbf{c}_k - \mathbf{x}_i),$$

expressing the total work done by the forces \mathbf{u}_{ki} from \mathbf{x}_i to \mathbf{c}_k , which is the potential energy of the centers configuration.

This mechanical model is a generalization to several centers of the Varignon Frame, [4, Section 1.3.4].

5. Strong duality

Theorem 2 is a **weak duality** theorem in the sense that any feasible solution \mathbf{U} of (D) gives a lower bound for the optimal value of (P), and conversely, any set of centers $\{\mathbf{c}_k\}$ for (P) gives an upper bound on the optimal value of (D). The next two theorems show that there is no duality gap between (P) and (D).

Theorem 3. *Given the data $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, and an optimal solution*

$$\{\mathbf{c}_1, \dots, \mathbf{c}_K\}, \quad \{p_{ki} : i \in \overline{1, N}, k \in \overline{1, K}\},$$

of the primal problem (P), there exist feasible dual variables \mathbf{U} such that

$$g(\mathbf{U}) = f(\mathbf{c}_1, \dots, \mathbf{c}_K). \tag{26}$$

Proof. We distinguish two cases.

Case 1. None of the centers $\{\mathbf{c}_k\}$ coincides with any of the data points $\{\mathbf{x}_i\}$.

For $k \in \overline{1, K}$ and $i \in \overline{1, N}$ define

$$\mathbf{u}_{ki} := \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{c}_k)} Q_k (\mathbf{x}_i - \mathbf{c}_k). \tag{27}$$

Then from (17a) and Theorem 1 it follows that,

$$\sum_{i=1}^N \mathbf{u}_{ki} = \mathbf{R}_k(\mathbf{c}_k) = \mathbf{0}, \text{ verifying (21).}$$

Rewriting (27) as

$$Q_k^{-1/2} \mathbf{u}_{ki} = \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{c}_k)} Q_k^{1/2} (\mathbf{x}_i - \mathbf{c}_k),$$

we get, for all k, i ,

$$\|Q_k^{-1/2} \mathbf{u}_{ki}\| = \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{c}_k)} \|\mathbf{x}_i - \mathbf{c}_k\|_{Q_k} = p_{ki}^2, \tag{28}$$

proving that the inequalities (22) hold as equations, and that $\{\mathbf{u}_{ki}\}$, defined by (27), are feasible.

From (27) and (4) and (5) it follows that:

$$\mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k) = p_{ki}^2 d_k(\mathbf{x}_i, \mathbf{c}_k), \tag{29}$$

and (26) follows from (24).

Case 2. A center coincides with one of the data points, say

$$\mathbf{c}_k = \mathbf{x}_j, \text{ for some } k \in \overline{1, K}, \quad j \in \overline{1, N}, \tag{30}$$

in which case (16) holds. Define

$$\mathbf{u}_{ki} := \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{x}_j)} Q_k (\mathbf{x}_i - \mathbf{x}_j), \text{ for } i \neq j, \tag{31a}$$

$$\mathbf{u}_{kj} := - \sum_{i \neq j} \mathbf{u}_{ki}. \tag{31b}$$

Then $\sum_i \mathbf{u}_{ki} = \mathbf{0}$ by definition, and $\|Q_k^{-1/2} \mathbf{u}_{ki}\| = p_{ki}^2$ for all $i \neq j$, as in (28). Next,

$$\mathbf{u}_{kj} = -\mathbf{R}_k^j \text{ by (17c), and therefore by (17b),}$$

$$\mathbf{R}_k(\mathbf{x}_j) = \mathbf{0} \text{ implies } p_{kj}^2 \geq \|Q_k^{-1/2} \mathbf{R}_k^j\| = \|Q_k^{-1/2} \mathbf{u}_{kj}\|,$$

proving that the variables \mathbf{U} defined by (31a) and (31b) are feasible.

Finally, we prove (26). As in Case 1 we have the equality (29) for all $i \neq j$. Also,

$$0 = \mathbf{u}_{mj} \cdot (\mathbf{x}_j - \mathbf{c}_m) \leq \|Q_m^{-1/2} \mathbf{u}_{mj}\| \|\mathbf{x}_j - \mathbf{c}_m\|_{Q_m} \leq p_{mj}^2 d_m(\mathbf{x}_j, \mathbf{c}_m) = 0$$

for any other cluster $m \in \overline{1, K}$, $m \neq k$, since $p_{mj} = 0$ (by (16)), and therefore $\mathbf{u}_{mj} = \mathbf{0}$.

The inequalities (25a)–(25b) therefore reduce to

$$\mathbf{u}_{kj} \cdot (\mathbf{x}_j - \mathbf{c}_k) \leq \|Q_k^{-1/2} \mathbf{u}_{kj}\| \|\mathbf{x}_j - \mathbf{c}_k\|_{Q_k} \leq p_{kj}^2 d_k(\mathbf{x}_j, \mathbf{c}_k),$$

which become trivial equalities, since by (30) all three terms are zero.

The inequalities (25a) and (25b) therefore hold as equalities, proving (26). \square

Now a converse of Theorem 3.

Theorem 4. Let \mathbf{U} be an optimal solution of the dual problem (D). Then there exist $\{\mathbf{c}_1, \dots, \mathbf{c}_K\}$ such that

$$g(\mathbf{U}) = f(\mathbf{c}_1, \dots, \mathbf{c}_K). \tag{26}$$

Proof. Writing the objective function $g(\mathbf{U})$ as in (24), the Lagrangian of (D) is

$$\sum_{k=1}^K \sum_{i=1}^N \mathbf{u}_{ki} \cdot (\mathbf{x}_i - \mathbf{c}_k) - \sum_{k=1}^K \sum_{i=1}^N t_{ki} (\|Q_k^{-1/2} \mathbf{u}_{ki}\| - p_{ki}^2)$$

with Lagrange multipliers $\{t_{ki}\}$. The Karush–Kuhn–Tucker necessary conditions for optimality are

$$(\mathbf{x}_i - \mathbf{c}_k) - t_{ki} Q_k^{-1} \frac{\mathbf{u}_{ki}}{\|Q_k^{-1/2} \mathbf{u}_{ki}\|} = \mathbf{0}, \tag{32a}$$

$$\sum_{i=1}^N \mathbf{u}_{ki} = \mathbf{0}, \tag{32b}$$

$$\|Q_k^{-1/2} \mathbf{u}_{ki}\| \leq p_{ki}^2, \tag{32c}$$

$$t_{ki} \geq 0, \tag{32d}$$

$$t_{ki} (\|Q_k^{-1/2} \mathbf{u}_{ki}\| - p_{ki}^2) = 0, \tag{32e}$$

for all $k \in \overline{1, K}$, $i \in \overline{1, N}$. Again we distinguish two cases.

Case 1. All $t_{ki} > 0$. Then (28) follows from (32e) for all k, i , and from (32a):

$$Q_k^{1/2} (\mathbf{x}_i - \mathbf{c}_k) = t_{ki} Q_k^{-1/2} \frac{\mathbf{u}_{ki}}{\|Q_k^{-1/2} \mathbf{u}_{ki}\|}.$$

Taking norms on both sides gives

$$d_k(\mathbf{x}_i, \mathbf{c}_k) = t_{ki}, \tag{33}$$

and by substituting (28) and (33) in (32a),

$$\mathbf{u}_{ki} := \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{c}_k)} Q_k (\mathbf{x}_i - \mathbf{c}_k),$$

and the equality (26) follows as in the proof of Theorem 3, Case 1.

Case 2. Some Lagrange multipliers are zero, say $t_{ki} = 0$. Then by (32a), $\mathbf{c}_k = \mathbf{x}_j$, and $t_{ki} > 0$ for $i \neq j$, by (32a) and (32d), and therefore, by (32e),

$$\|Q_k^{-1/2} \mathbf{u}_{ki}\| = p_{ki}^2, \text{ for all } i \neq j.$$

From $\mathbf{c}_k = \mathbf{x}_j$ and (32a) it follows that:

$$Q_k^{1/2} (\mathbf{x}_i - \mathbf{x}_j) = t_{ki} Q_k^{-1/2} \frac{\mathbf{u}_{ki}}{\|Q_k^{-1/2} \mathbf{u}_{ki}\|}$$

and by taking norms,

$$d_k(\mathbf{x}_i, \mathbf{x}_j) = t_{ki}, \text{ for all } i \neq j.$$

Substituting t_{ki} and $\|Q_k^{-1/2} \mathbf{u}_{ki}\|$ in (32a) gives,

$$\mathbf{u}_{ki} := \frac{p_{ki}^2}{d_k(\mathbf{x}_i, \mathbf{x}_j)} Q_k (\mathbf{x}_i - \mathbf{x}_j), \text{ for } i \neq j,$$

and from (32b),

$$\mathbf{u}_{kj} := - \sum_{i \neq j} \mathbf{u}_{ki},$$

reproducing (31a) and (31b), and equality in (26) follows as in the proof of Theorem 3, Case 2. \square

6. Discussion

- (a) The practical applicability of Theorems 2, 3, 4 above is limited by the fact that the matrices $\{Q_1, \dots, Q_K\}$ modelling the geometry of the clusters are not known a priori, and may require a solution of the primal problem (P). However, useful bounds on the optimal value of (P) can be found by taking Euclidean distances, i.e., approximating the matrices $\{Q_k\}$ by the identity matrix.
- (b) For other results on duality in multi-facility location problems see [9,12] and their references.
- (c) A matrix analog of the harmonic mean is the **parallel sum** of Anderson and Duffin [1], defined for matrices $A, B \in \mathbb{C}^{n \times n}$ by

$$A:B = A(A+B)^\dagger B, \quad (34)$$

where † denotes the Moore–Penrose inverse, see also [6]. The parallel sum of 3 or more matrices is defined inductively.

The JDF (9) can be written as a parallel sum. Consider, without loss of generality, the case of 2 clusters. Let $D_1 = \text{diag}(d_1(\mathbf{x}_i))$, $D_2 = \text{diag}(d_2(\mathbf{x}_i))$ be $N \times N$ diagonal matrices, with diagonal elements the distances of the data points \mathbf{x}_i to the centers $\mathbf{c}_1, \mathbf{c}_2$. Then the diagonal elements of the parallel sum $D_1 : D_2$ are the values $\{J(\mathbf{x}_i)\}$ of (9) with $K = 2$.

Extensions to more general (not diagonal) distance matrices may be of interest.

(d) The problem (P) can be written as

$$\mathbf{1} \cdot (D_1 : D_2) \mathbf{1} = \min_{\substack{\mathbf{p}_1, \mathbf{p}_2 \\ \mathbf{p}_1 + \mathbf{p}_2 = \mathbf{1}}} \{\mathbf{p}_1 \cdot D_1 \mathbf{p}_1 + \mathbf{p}_2 \cdot D_2 \mathbf{p}_2\}, \quad (35)$$

where $\mathbf{p}_1 = (p_1(\mathbf{x}_i))$, $\mathbf{p}_2 = (p_2(\mathbf{x}_i))$ are the vectors of probabilities, and $\mathbf{1}$ is a vector of ones.

Eq. (6) admits a simple physical analog. Consider an electrical circuit consisting of resistances (resistance matrices) D_1, D_2 connected in parallel. The resistance of the circuit is the parallel sum $D_1 : D_2$, see [1]. The left side of (6) is the energy dissipated by a vector current $\mathbf{1}$ through this circuit.

The incoming current $\mathbf{1}$ splits to currents \mathbf{p}_1 and \mathbf{p}_2 through the resistances D_1 and D_2 , respectively. According to Kelvin's principle, these currents minimize the energy dissipated in the circuit, given by the right side of (6). For related results see [10].

References

- [1] W.N. Anderson Jr., R.J. Duffin, Series and parallel addition of matrices, *J. Math. Anal. Appl.*, 26 (1969) 576–594.
- [2] M. Arav, Contour approximation of data and the harmonic mean, *J. Math. Inequalities* 2 (2008) 161–167.
- [3] A. Ben-Israel, C. Iyigun, Probabilistic distance clustering, *J. Classification* 25 (2008) 5–26.
- [4] Z. Deznier, K. Klamroth, A. Schöbel, G.O. Wesolowsky, The Weber problem, in: Z. Deznier, H.W. Hamacher (Eds.), *Facility Location: Applications and Theory*, Springer, 2002 (Chapter 1).
- [5] C. Iyigun, A. Ben-Israel, Probabilistic distance clustering adjusted for cluster size, *Probab. Engrg. Info. Sci.* 22 (2008) 1–19.
- [6] F. Kubo, T. Ando, Means of positive linear operators, *Math. Ann.* 246 (1980) 205–224.
- [7] H.W. Kuhn, On a pair of dual nonlinear programs, in: J. Abadie (Ed.), *Methods of Nonlinear Programming*, North-Holland, Amsterdam, 1967, pp. 38–54.
- [8] H.W. Kuhn, A note on Fermat's problem, *Math. Program.* 4 (1973) 98–107.
- [9] R.F. Love, H. Juel, Properties and solution methods for large location–allocation problems, *J. Oper. Res. Soc.* 33 (1982) 443–452.
- [10] T.D. Morley, Parallel summation, Maxwells principle and the infimum of projections, *J. Math. Anal. Appl.* 70 (1979) 33–41.
- [11] M. Teboulle, A unified continuous optimization framework for center-based clustering methods, *J. Mach. Learn.* 8 (2007) 65–102.
- [12] H. Üster, R.F. Love, Duality in constrained multi-facility location models, *Naval Res. Logist. Quart.* 49 (2002) 410–421.
- [13] E. Weiszfeld, Sur le point par lequel la somme des distances de n points donnés est minimum, *Tohoku Math. J.* 43 (1937) 355–386.