

# DISTANCE BASED PROBABILISTIC CLUSTERING OF DATA

CEM IYIGUN AND ADI BEN-ISRAEL

The problem is to partition a given data set  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subset \mathbb{R}^n$ , into clusters  $\{\mathcal{C}_1, \dots, \mathcal{C}_m\}$ , where the number  $m$  of clusters is either given, or to be determined by an optimality criterion.

Clusters consist of similar points, and are themselves dissimilar, where similarity is in a sense of a distance  $d(\cdot, \cdot)$  on  $\mathbb{R}^n$ .

With a cluster  $\mathcal{C}_i$ , we associate a **center**  $\mathbf{c}_i$ , and for any data point  $\mathbf{x} \in \mathcal{D}$  we then compute:

- a **distance**  $d(\mathbf{x}, \mathbf{c}_i)$ , denoted by  $d_i(\mathbf{x})$ , and
- a **probability**  $p_i(\mathbf{x})$  of membership in  $\mathcal{C}_i$ .

We assume throughout that for all  $\mathbf{x}$ ,

$$p_i(\mathbf{x}) d_i(\mathbf{x}) = \text{constant, depending on } \mathbf{x}, \quad \text{for } i = 1, \dots, m, \quad (1)$$

making membership in nearby clusters more probable. Since probabilities add to 1, assumption (1) implies

$$p_i(\mathbf{x}) = \frac{\prod_{j \neq i} d_j(\mathbf{x})}{\sum_{k=1}^m \prod_{j \neq k} d_j(\mathbf{x})}, \quad i = 1, \dots, m, \quad (2)$$

in particular, for  $m = 2$ ,

$$p_1(\mathbf{x}) = \frac{d_2(\mathbf{x})}{d_1(\mathbf{x}) + d_2(\mathbf{x})}, \quad p_2(\mathbf{x}) = \frac{d_1(\mathbf{x})}{d_1(\mathbf{x}) + d_2(\mathbf{x})}. \quad (3)$$

We present a new clustering algorithm that iterates on centers, distances and probabilities, compare it with existing methods, and illustrate its advantages.

*Keywords:* Distance based clustering, probabilistic clustering.

*E-mail address:* {iyigun,bisrael}@rutcor.rutgers.edu

RUTCOR—RUTGERS CENTER FOR OPERATIONS RESEARCH, RUTGERS UNIVERSITY, 640 BARTHOLOMEW RD., PISCATAWAY, NJ 08854-8003, USA