

A NOVEL DATA REPRESENTATION BASED ON DISSIMILARITY INCREMENTS

Helena Aidos and Ana Fred

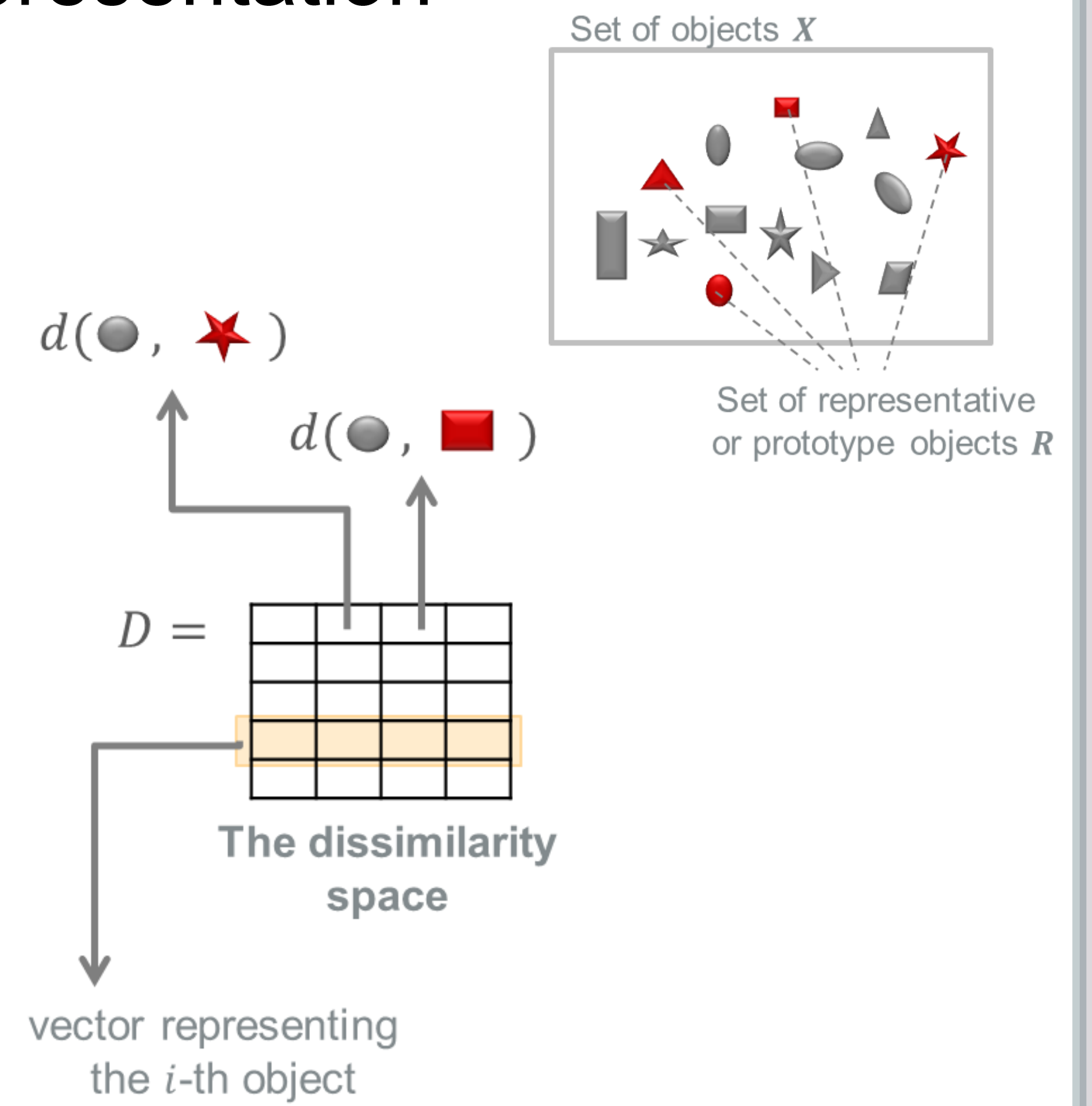
Instituto de Telecomunicações
Instituto Superior Técnico, Universidade de Lisboa, Portugal
{haidos, afred}@lx.it.pt

Motivation

- Typically, objects are represented by a set of features, which should characterize the objects and be relevant to discriminate among the classes.
- Problem:** difficult to obtain a complete description of objects:
 - forces an overlap of the classes
 - leads to an inefficient learning process.
- Solution:** Use a dissimilarity representation, which is based on comparisons between pairs of objects:
 - Solves the problem of class overlap, since only identical objects have a dissimilarity of zero.

Dissimilarity representation

- $X = \{x_1, \dots, x_n\}$ set of objects
- $R = \{e_1, \dots, e_r\}$ set of representative or prototype objects, such that $R \subseteq X$
- Each object x_i is described by a r -dimensional dissimilarity vector
 $D(x_i, R) = [d(x_i, e_1) \dots d(x_i, e_r)]$
where $d(\cdot, \cdot)$ is a dissimilarity measure
 - $D(x_i, R)$ is a row of the $n \times r$ dissimilarity matrix D , the **dissimilarity space**
- Define a vector space Y by $Y = D$, where the i -th object is represented by the dissimilarity vector of the D_{ij} values.



PROPOSAL: A novel dissimilarity representation of data, based on a second-order dissimilarity measure.

Second-order dissimilarity measure: the dissimilarity increments

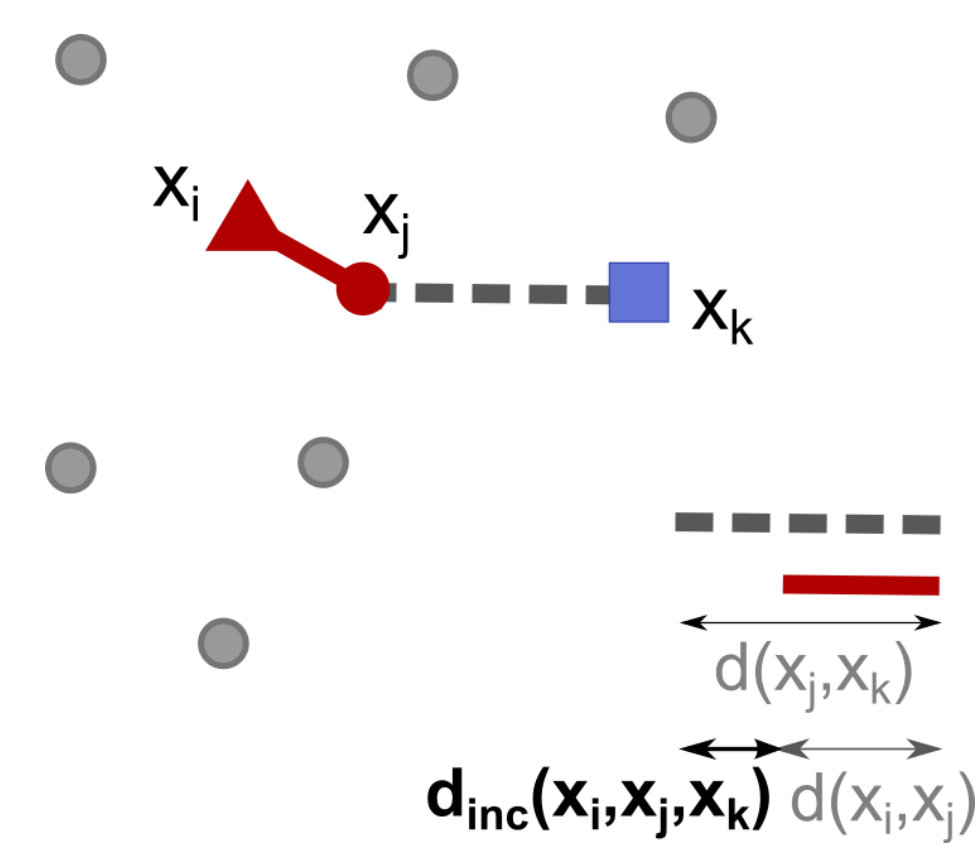
Given some dissimilarity measure, $d(\cdot, \cdot)$, between patterns,

(x_i, x_j, x_k) – triplet of nearest neighbor

- x_j is the nearest neighbor of x_i
- x_k is the nearest neighbor of x_j (different from x_i)

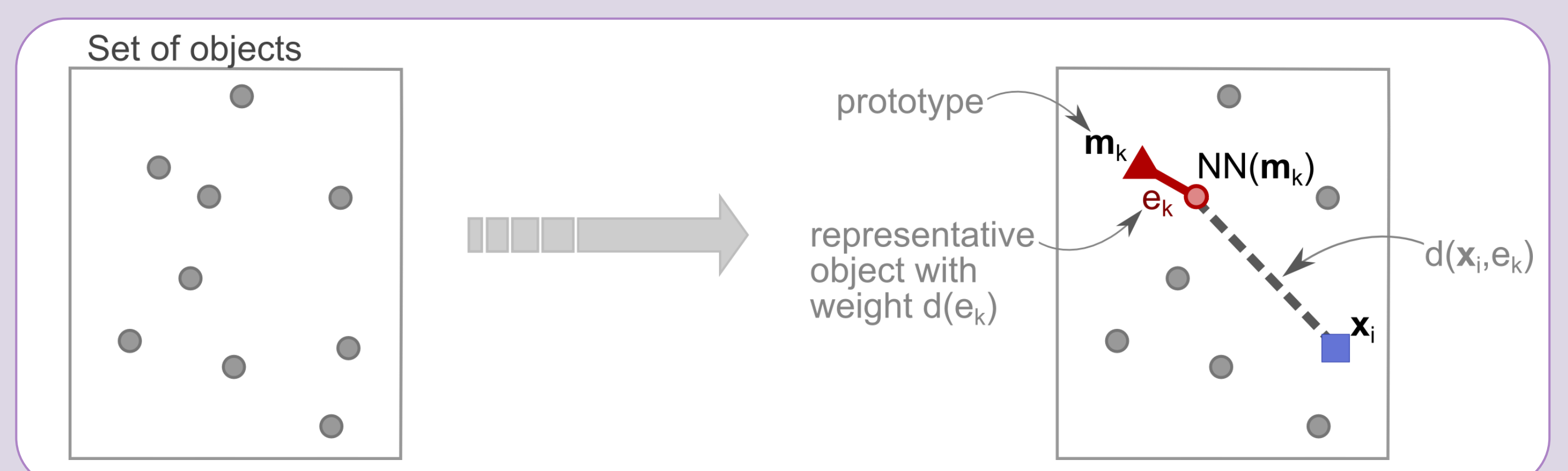
The **dissimilarity increments** between neighboring patterns is defined as

$$d_{inc}(x_i, x_j, x_k) = |d(x_i, x_j) - d(x_j, x_k)|$$



Dissimilarity increments space

- $R = \{e_1, \dots, e_r\}$ set of prototype objects, with e_j an edge between a prototype m_j and its nearest neighbor $x_{m_j} = NN(m_j)$
- $d(e_j) = d(m_j, x_{m_j})$ weight of edge e_j
- Distance between any object x_i and the representative object e_j is
 $d(x_i, e_j) = \min \{d(x_i, m_j), d(x_i, x_{m_j})\}$



The (i, j) -th element of the Dinc space is defined as
 $D(x_i, e_j) = |d(x_i, e_j) - d(e_j)|$

Characterization

- Dissimilarity spaces have **higher discriminant power** of features in separating the classes.
- Dissimilarity spaces have **less overlap between the classes**, which may facilitate the learner to separate the samples of different classes.
- Even if the classes are more separable, they are nonlinearly separable by 1-NN classifier.

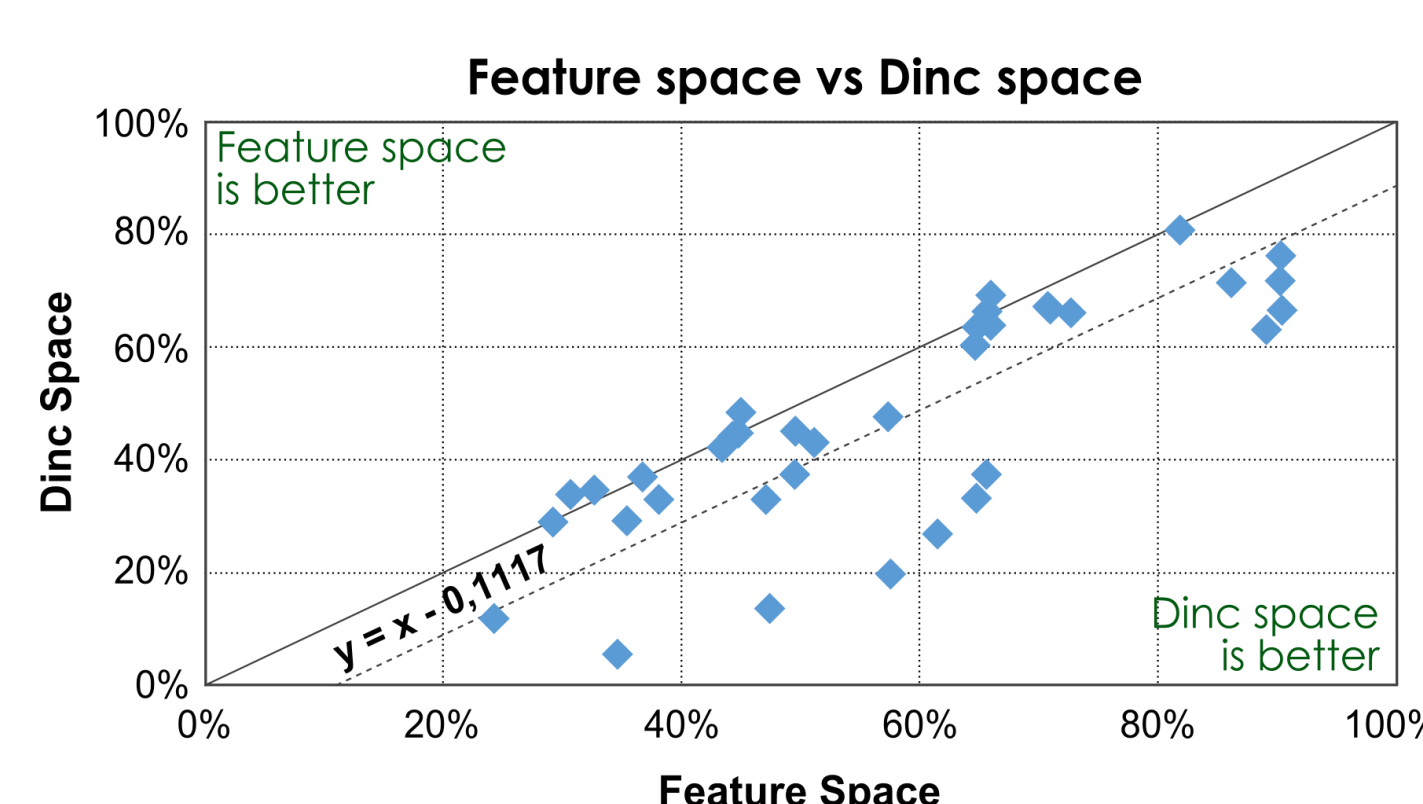
Assume that $R = X$, meaning that all objects of X are used as prototypes.

Euclidean dissimilarity space

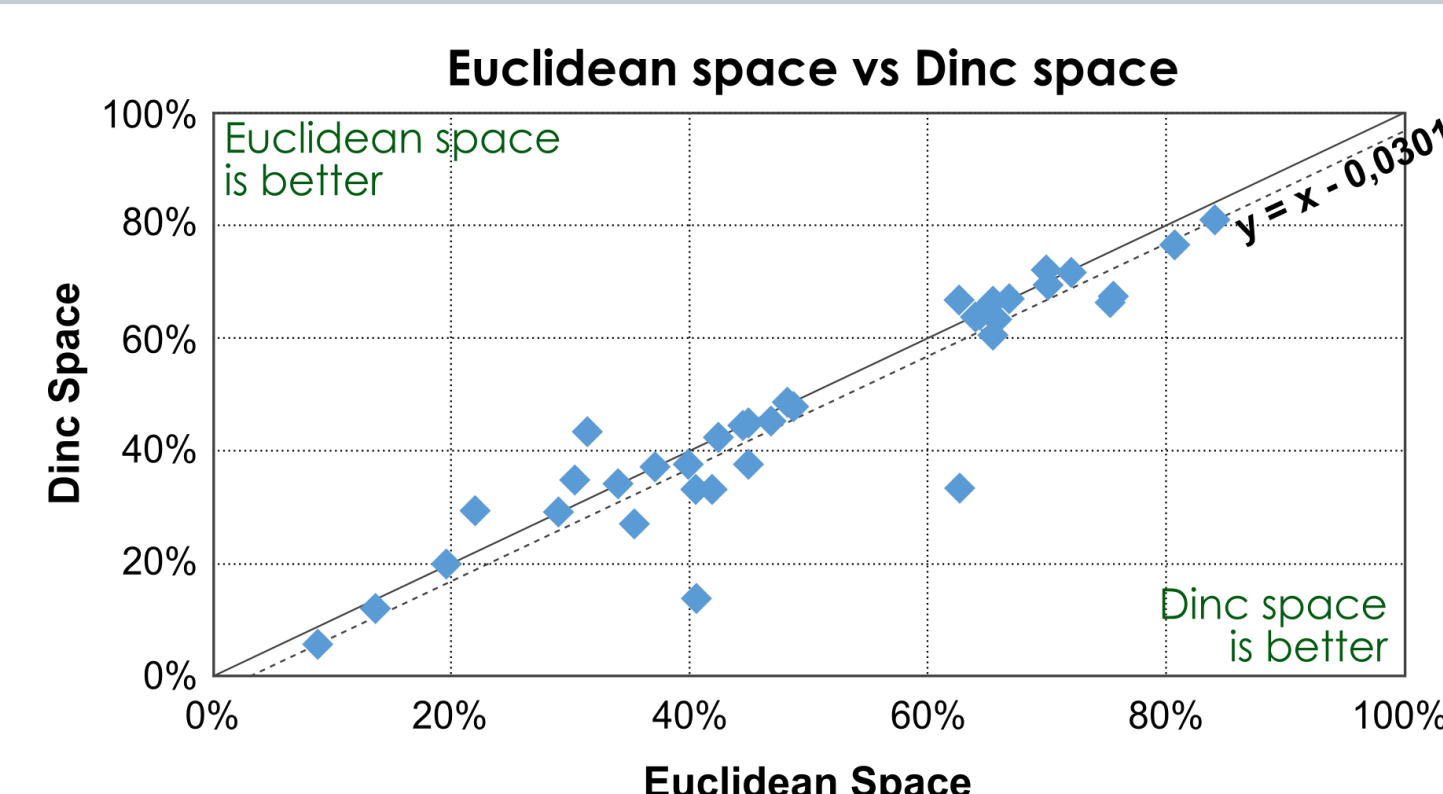
Each element, D_{ij} , of the dissimilarity matrix D , is the Euclidean distance between i -th and j -th objects.

Datasets: 36 real-world datasets from the UCI Machine Learning repository
Evaluation: Error rates of median-link, when the true number of clusters is known

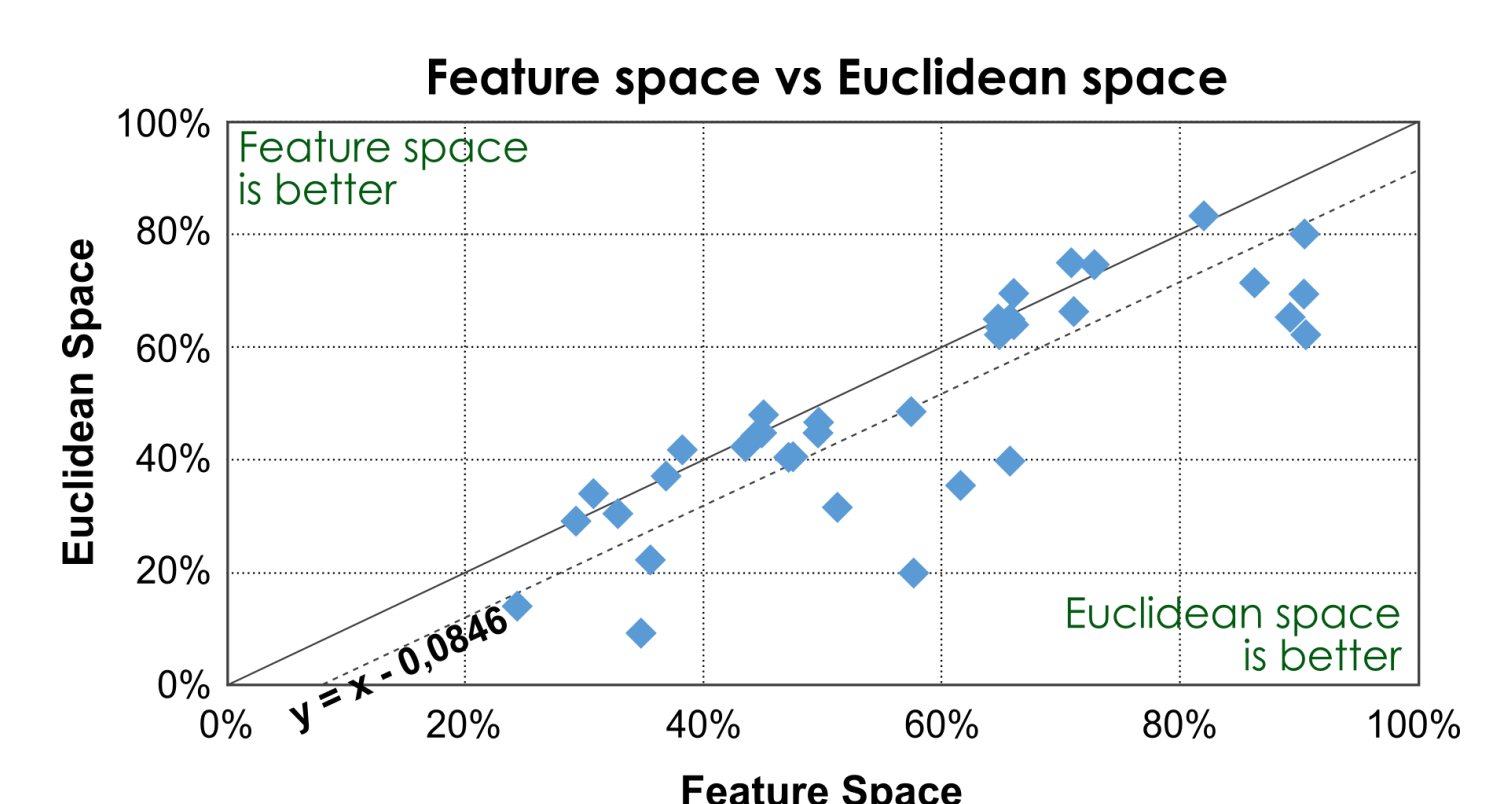
Experimental results



- 28 datasets for Dinc space
 - Best on average 13.6% than Feature space
- 6 datasets for Feature space
 - Best on average 2.2% than Dinc space



- 8 datasets for Euclidean space
 - Best on average 4.0% than Dinc space
- 18 datasets for Dinc space
 - Best on average 7.1% than Euclidean space



- 25 datasets for Euclidean space
 - Best on average 11.8% than Feature space
- 9 datasets for Feature space
 - Best on average 2.6% than Euclidean space