

A NOVEL DATA REPRESENTATION BASED ON DISSIMILARITY INCREMENTS

Helena Aidos and Ana Fred

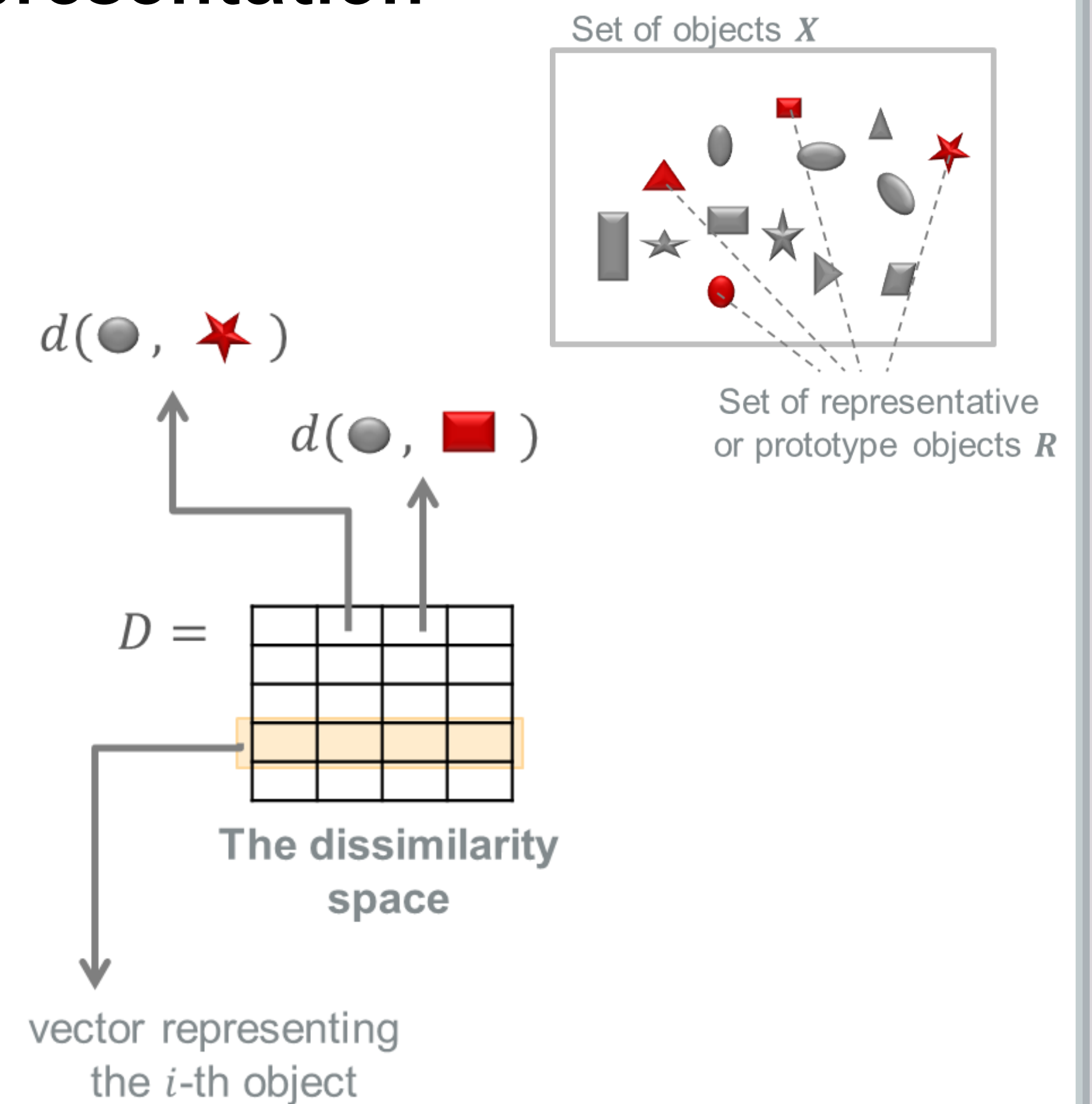
Instituto de Telecomunicações
Instituto Superior Técnico, Universidade de Lisboa, Portugal
{haidos, afred}@lx.it.pt

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

- set of objects
- set of representative or prototype objects, such that
- Each object is described by a d -dimensional dissimilarity vector (d_1, d_2, \dots, d_d) where d_{ij} is a dissimilarity measure
- d_{ij} is a row of the dissimilarity matrix D , the **dissimilarity space**
- Define a vector space by D , where the i -th object is represented by the dissimilarity vector of the i -th values.



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

Second-order dissimilarity measure: the dissimilarity increments

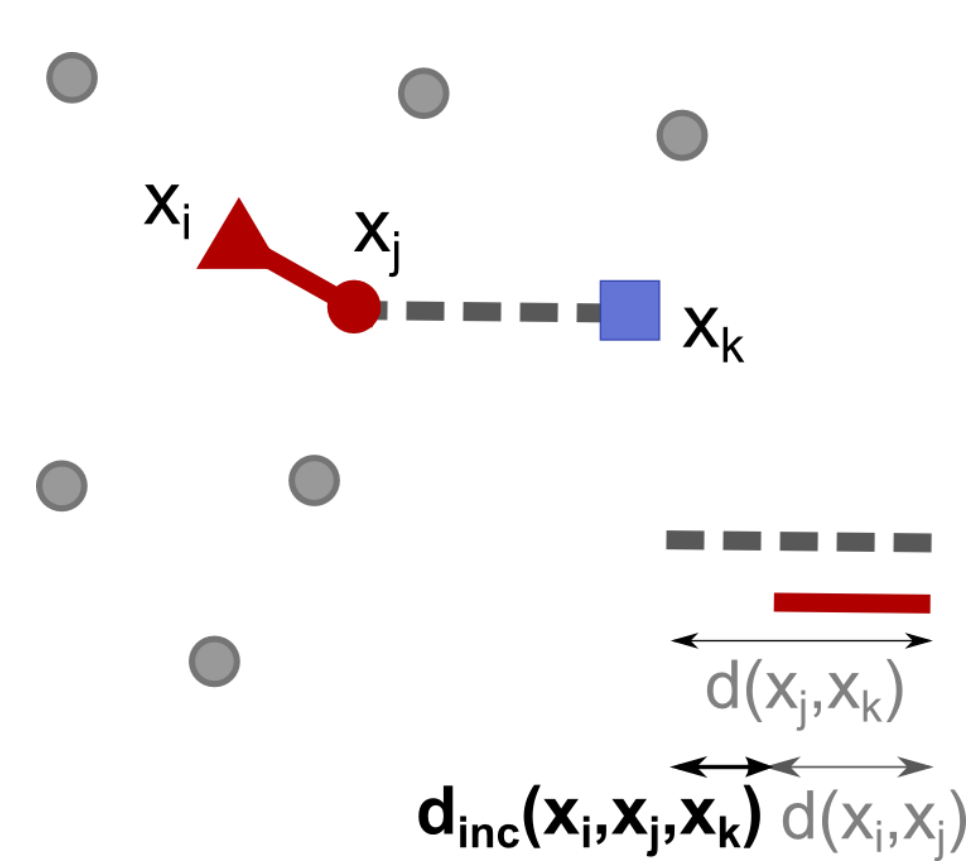
Given some dissimilarity measure, $d(x_i, x_j)$, between patterns,

(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)

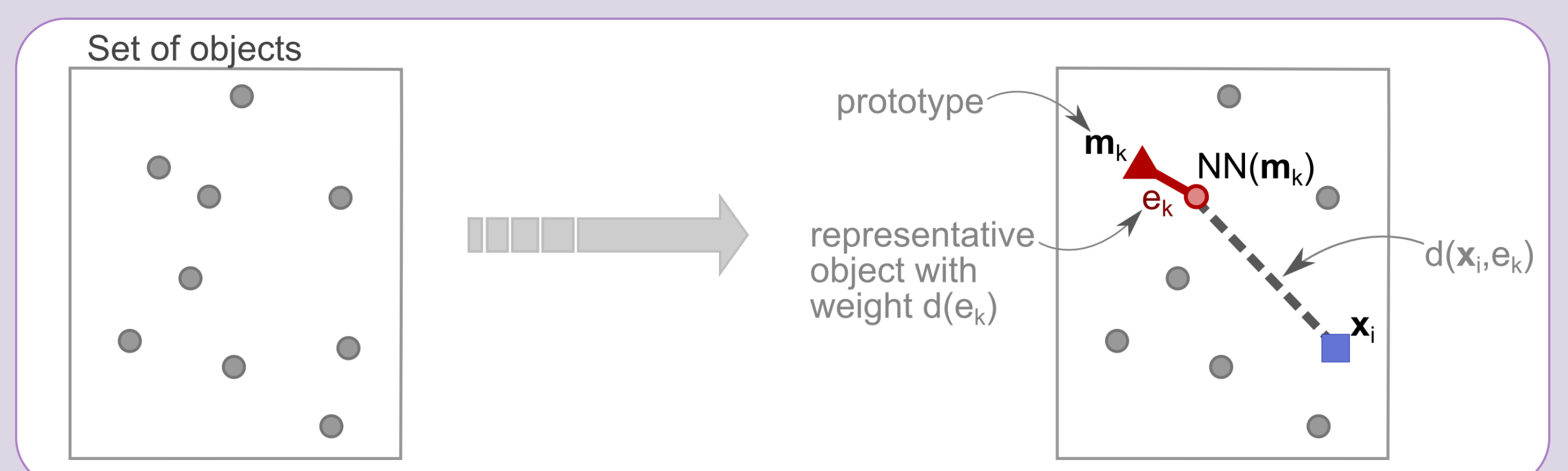
The **dissimilarity increments** between neighboring patterns is defined as

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



Dissimilarity increments space

- set of prototype objects, with an edge between a prototype and its nearest neighbor
- (e_k) weight of edge
- Distance between any object x_i and the representative object e_k is $d(x_i, e_k)$



The i -th element of the Dinc space is defined as $d(x_i, e_k)$

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.

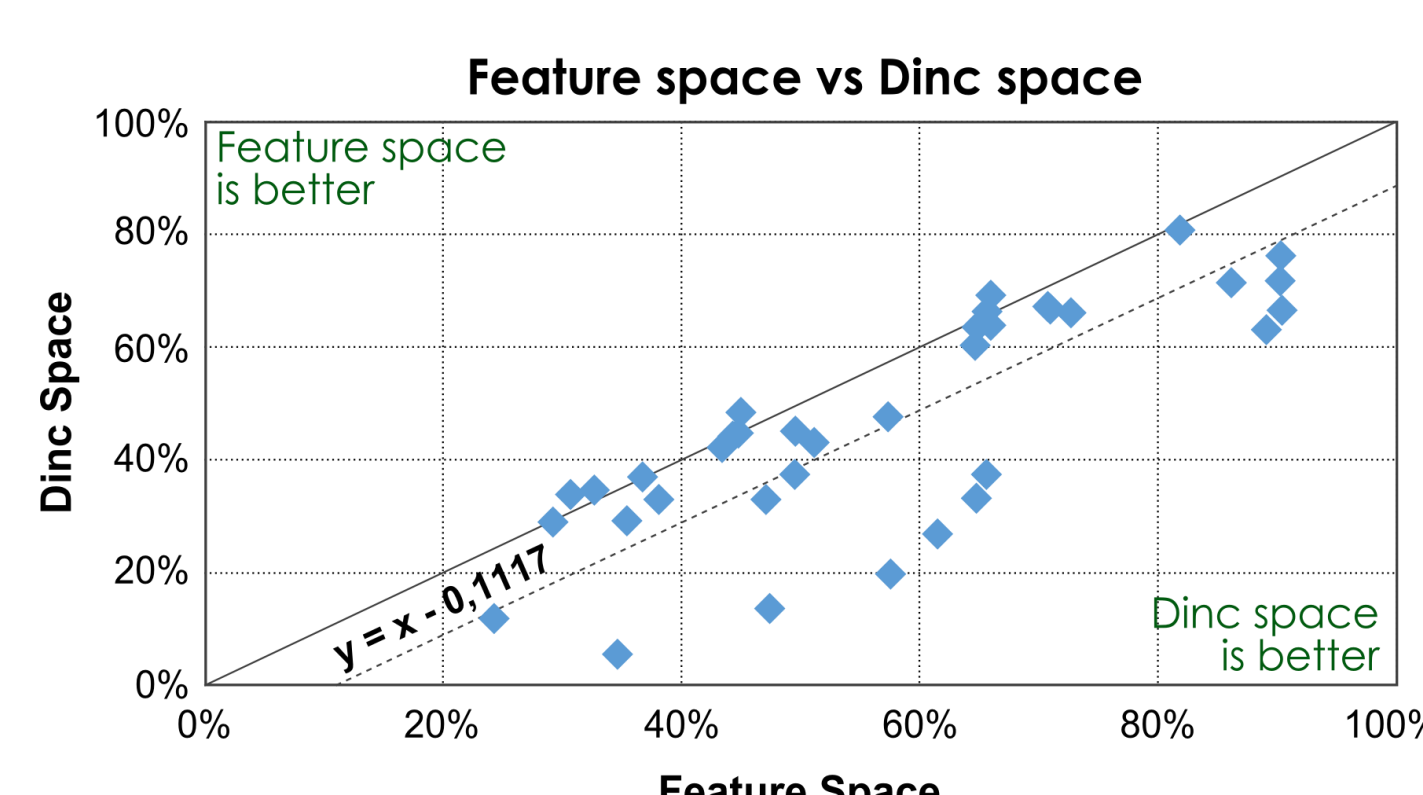
Assume that $R = X$, meaning that all objects of X are used as prototypes.

Euclidean dissimilarity space

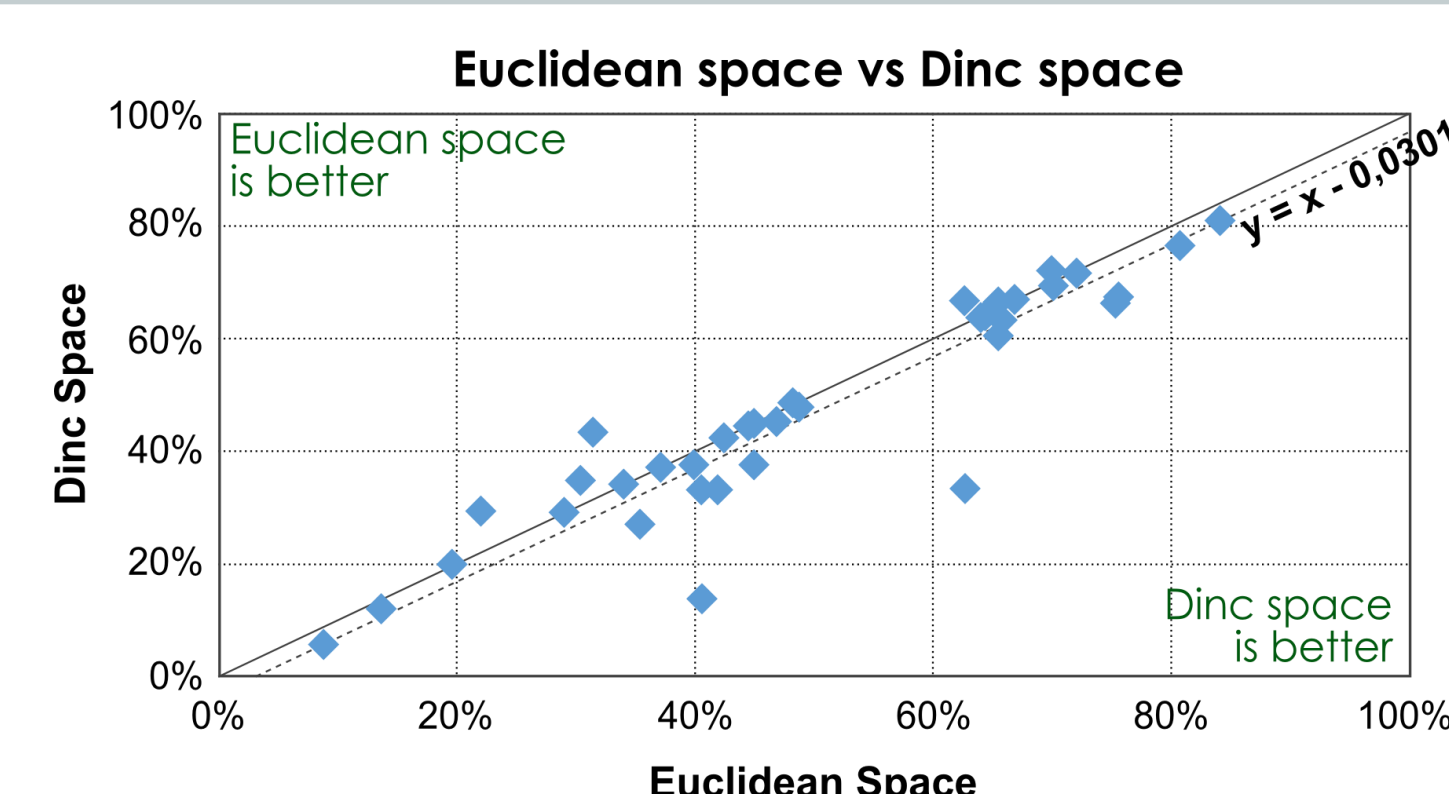
Each element, d_{ij} , of the dissimilarity matrix D , is the Euclidean distance between i -th and j -th objects.

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

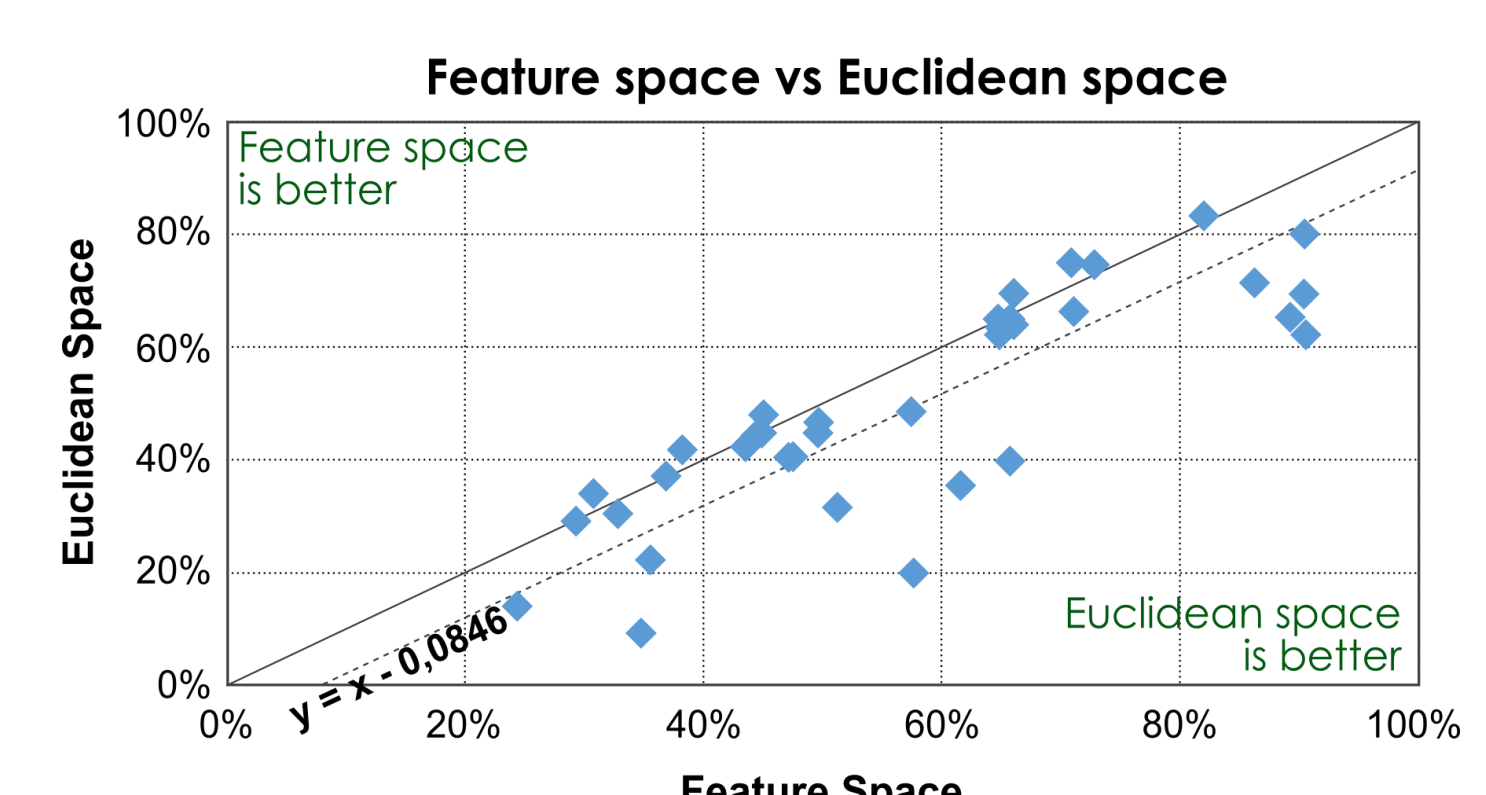
Experimental results



- 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