



NOVA
IMS
Information
Management
School

Data Mining

Semi-supervised classification

Nearest Neighbors

24/11/2021
NOVA-IMS
Fernando Lucas Bação
bação@isegi.unl.pt
<http://www.isegi.unl.pt/fbação>
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

1

NOVA
IMS
Information
Management
School

AGENDA

- Cluster analysis
 - Semi-supervised classification
 - Nearest Neighbors

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

2

NOVA
IMS
Information Management School

Semi-supervised classification

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

UNIGES, A3ES, eSchools, eduniversal

3

NOVA
IMS
Information Management School

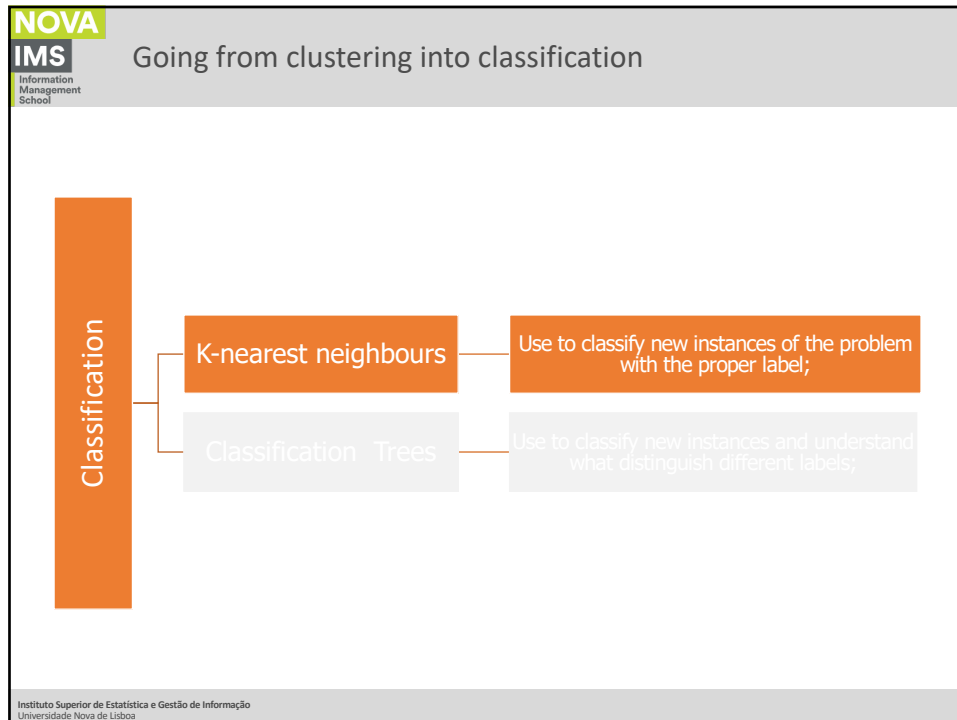
Going from clustering into classification

```
graph LR; A[Unlabeled Data] --> B[Clustering - Labels]; B --> C[Classification - KNN and Trees]
```

Unlabeled Data Clustering - Labels Classification – KNN and Trees

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

4



5

NOVA
IMS
Information Management School

k-nearest neighbors (k-NN)

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

6

NOVA

IMS
Information
Management
School

k -nearest neighbors

- Instance based classification:
 - Simplest form of learning;
 - Training instances are searched for instances that most closely resembles new instance;
 - The instances themselves represent the knowledge;
 - Also called instance-based learning
 - Similarity function defines what's "learned"

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

7

NOVA

IMS
Information
Management
School

k -nearest neighbors

- Requires three things:
 1. The set of stored records (with labels)
 2. A distance metric to compute distance between records (can use Euclidean distance)
 3. The value of k , the number of nearest neighbors to retrieve

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

8

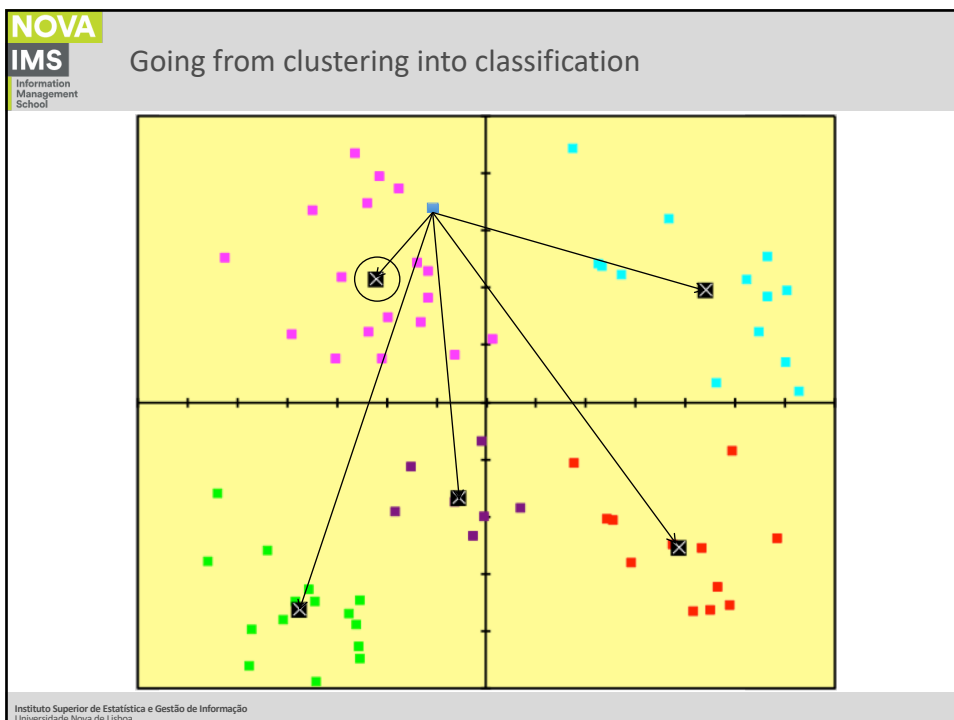
NOVA
IMS
Information Management School

k -nearest neighbors

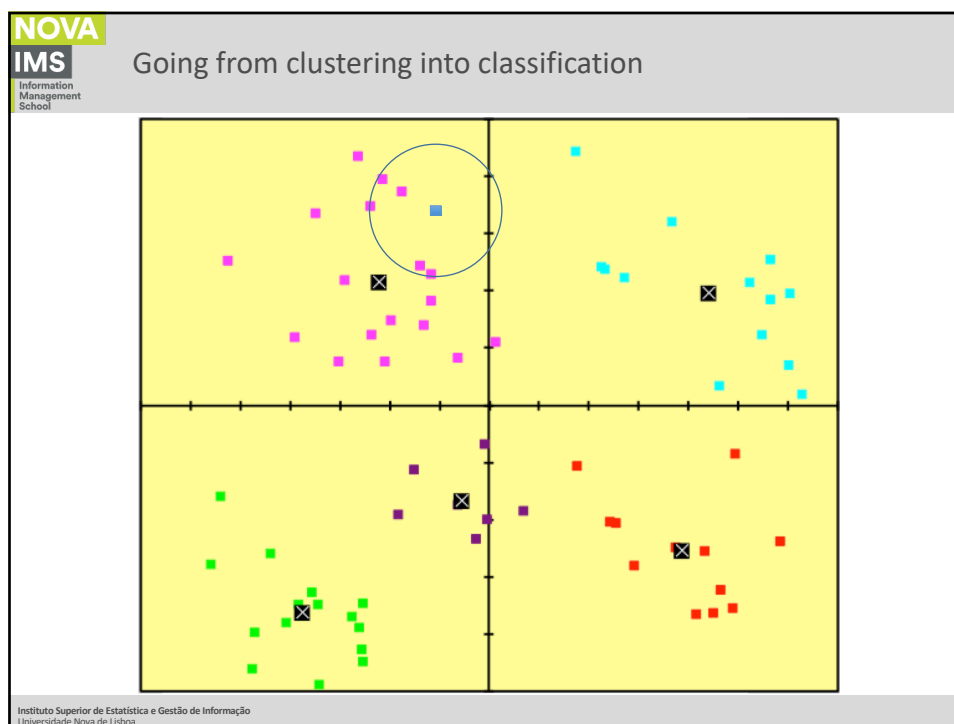
- To classify an unknown record:
 1. Compute distance to other training records
 2. Identify k nearest neighbors
 3. Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

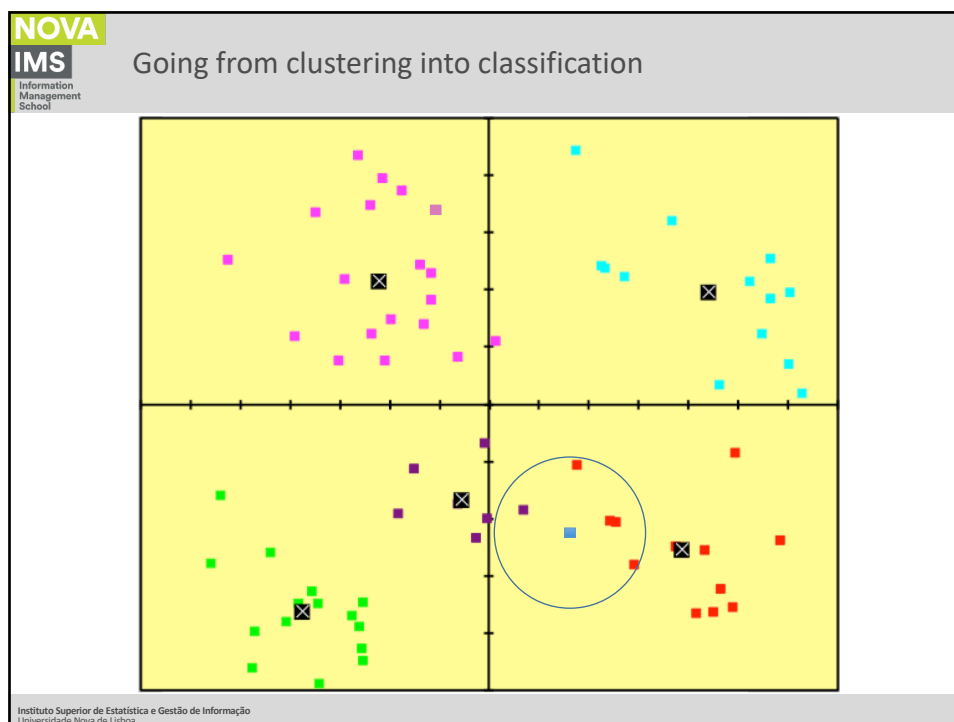
9



10



11



12

NOVA
IMS
Information Management School

k-nearest neighbors

- Compute distance between two points:
 - Euclidean distance

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$
- Determine the class from nearest neighbor list
 - Take the majority vote of class labels among the *k*-nearest neighbors
 - Weigh the vote according to distance
 - weight factor, $w = 1/d^2$

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

13

NOVA
IMS
Information Management School

k-nearest neighbors

- *k*-nn frontiers (and the number *k*):
 - Large *k*
 - Smooth frontiers
 - Unable to detect small variations
 - Small *k*
 - Very sensitive to outliers
 - Crisp frontiers

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

14

