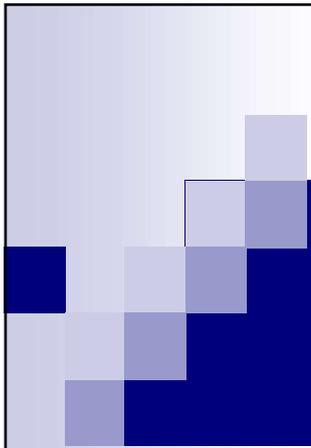


# Sistemas de Apoio à Decisão– Clustering

V 1.0, V.Lobo, EN, 2018



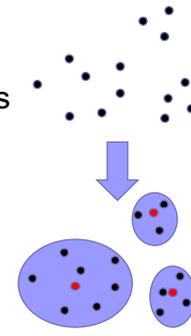
## Clustering (agrupamento)

Victor Lobo

1

## Clustering

- Objectivo fundamental
  - Definir agrupamentos de dados que têm **algo em comum ou parecido**
    - Sumarizar
    - Detectar grandes grupos / detectar outliers
    - Facilitar a gestão de múltiplas entidades
- Também conhecido como:
  - Técnicas de Agrupamento
  - Na comunidade de estatística: Classificação (dividir em classes)



2



## Tipos de técnicas

### ■ Hierárquicas

- Sub-divisões cada vez mais detalhadas
  - “Divisivas” – Começa no todo, e vai dividindo
  - “Aglomerativas” – Vai juntando grupos cada vez maiores
- *Dendogramas*, métodos de Ward, etc

### ■ Partições

- Divide o espaço em blocos sem relações hierárquicas
- *K-Médias*, Fuzzy C-mean, DBSCAN, SOM, etc

### ■ Outras (por densidade, por grelha, etc)

5

## Algoritmo “k-médias” de Lloyd

### ■ Algoritmo de Lloyd

1. Escolher aleatoriamente  $k$  centroides  $\mu_k$
2. Para cada exemplo  $x_j$ , encontrar o centroide  $\mu_k$  mais próximo, e atribuir-lhe a classe  $C_k$ .
3. Recalcular os centroides de cada classe
$$\mu_k = \sum_{x_j \in C_k} x_j / n_k$$
4. Voltar a 2, até que não haja alterações em  $\mu_k$

### ■ Exemplo

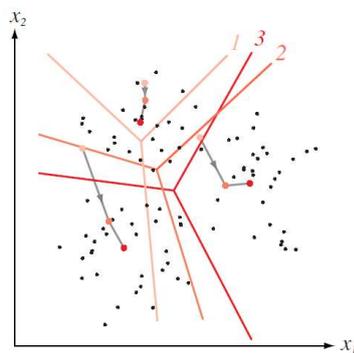
6

## Variantes

- Algoritmo de McQueen
  - actualizações incrementais
- K-medoids
  - Os centroids têm que ser pontos que existam nos dados originais
- Fuzzy c-Means
  - As pertenças aos centroids são funções fuzzy

7

## Exemplo



- 1) Para os centroides originais (a vermelho claro), as fronteiras são apresentadas em 1. Se calcularmos a média desses pontos, obtemos os centroides de 2.
- 2) Para o 2º conjunto de centroides (a vermelho médio), as fronteiras são apresentadas em 2. Se calcularmos a média desses pontos, obtemos os centroides de 3.
- 3) Para o 3º conjunto de centroides (a vermelho vivo), as fronteiras são apresentadas em 3. Ao calcularmos os novos centroides, as fronteiras não mudam.

8

8

## Classificadores hierarquicos

- Critérios para associar dados
  - Distância a um do grupo
  - Distância à média do grupo
  - Juntar sempre 2 a 2 grupos equivalentes
- Para definir os clusters
  - Escolher um nível de corte
- Exemplo para o Iris-Dataset

