

# MSC ENGENHARIA E CIÊNCIA DE DADOS

## MÉTODOS ESTATÍSTICOS EM DATA MINING

---

### Relatório 2

---

**Authors:**

Diogo Vilela (96193)

Pedro Rodrigues (96301)

José Pedro Antunes (96260)

Sebastião Caldas (96321)

Nuno Marques (95758)

[diogo.pimpao.vilela@tecnico.ulisboa.pt](mailto:diogo.pimpao.vilela@tecnico.ulisboa.pt)  
[pedro.maria.rodrigues@tecnico.ulisboa.pt](mailto:pedro.maria.rodrigues@tecnico.ulisboa.pt)  
[jose.pedro.m.c.a@tecnico.ulisboa.pt](mailto:jose.pedro.m.c.a@tecnico.ulisboa.pt)  
[sebastiao.caldas@tecnico.ulisboa.pt](mailto:sebastiao.caldas@tecnico.ulisboa.pt)  
[nuno.figueiredo.marques@tecnico.ulisboa.pt](mailto:nuno.figueiredo.marques@tecnico.ulisboa.pt)

Grupo 04

2022/2023 – 1º Semestre, P1

## CONTENTS

<b>I</b>	<b>Introdução</b>	<b>3</b>
<b>II</b>	<b>Objetivo</b>	<b>3</b>
<b>III</b>	<b>Métodos de Clustering</b>	<b>3</b>
III-A	Métodos Aglomerativos . . . . .	3
III-A1	Distância de Hamming . . . . .	3
III-A2	Distância de Gower . . . . .	3
III-B	Métodos de Partição . . . . .	4
III-B1	K-means . . . . .	4
III-C	DBSCAN . . . . .	4
<b>IV</b>	<b>Conclusão</b>	<b>4</b>
	<b>References</b>	<b>4</b>

## I. INTRODUÇÃO

## II. OBJETIVO

## III. MÉTODOS DE CLUSTERING

### A. Métodos Aglomerativos

Os métodos de *clustering* aglomerativos baseiam-se na fusão recursiva de *clusters* em cada nível hierárquico de acordo com uma dada métrica ou medida. Inicialmente, as observações são os seus próprios *clusters*, que se vão unindo, sucessivamente, em cada iteração do algoritmo, fazendo com que cada nível tenha menos um *cluster* do que aquele que lhe precede. O algoritmo termina, naturalmente, com um único *cluster*, ao qual todas as observações originais pertencem.

Evidentemente, estes métodos dependem da escolha da relação de *dissemelhança* entre objetos e do critério de fusão de *clusters*; prendendo-se, sobretudo, com a natureza dos dados em questão. Visto que os nossos dados são maioritariamente categóricos binários, e após terem sido retiradas as variáveis contínuas, foram abordadas duas estratégias: aplicar *one-hot encoding* aos dados, considerando como métrica a distância de *hamming* e; não aplicar nenhuma transformação e considerar a distância de *gower*.

Apesar de se terem formulado duas estratégias distintas, a metodologia permaneceu igual. Efetivamente, consideram-se os métodos de agrupamento *complete-linkage*, *single-linkage* e, *average-linkage*, escolhendo-se o número de *clusters* de modo a maximizar o *silhouette coefficient*. Seguidamente, em função do método que apresenta ter *clusters* mais distintos, comparou-se os valores da partição com os valores da variável resposta, obtendo-se as métricas adequadas.

1) *Distância de Hamming*: Em seguida apresentam-se os resultados para a primeira estratégia de *clustering* aglomerativo

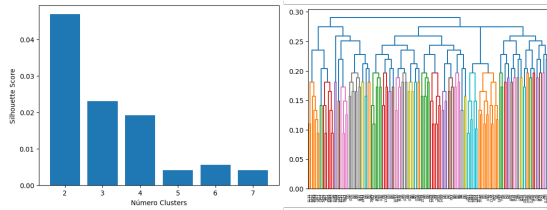


Fig. 1. *Silhouette Scores* e dendrograma (truncado) para *complete-linkage*.

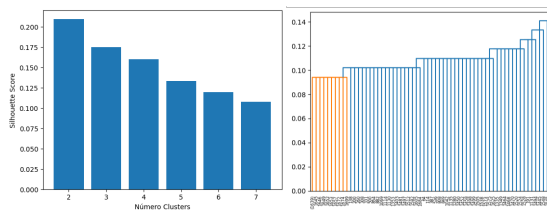


Fig. 2. *Silhouette Scores* e dendrograma (truncado) para *single-linkage*.

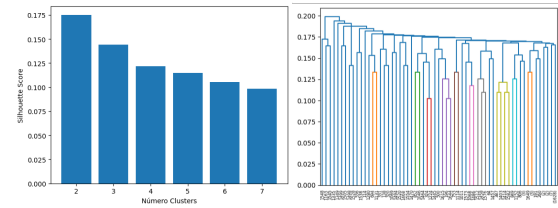


Fig. 3. *Silhouette Scores* e dendrograma (truncado) para *average-linkage*.

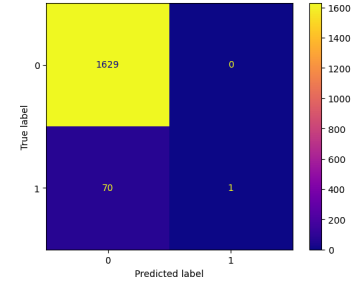


Fig. 4. Matriz de confusão para as *true labels* vs *clusters*.

	Random Index	Adjusted Random Index	Adjusted Mutual Index
Score	0.921	0.026	0.020

TABLE I  
MÉTRICAS PARA DISTÂNCIA DE HAMMING

2) *Distância de Gower*: Em seguida apresentam-se os resultados para a primeira estratégia de *clustering* aglomerativo

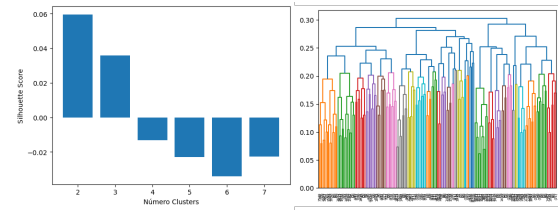


Fig. 5. *Silhouette Scores* e dendrograma (truncado) para *complete-linkage*.

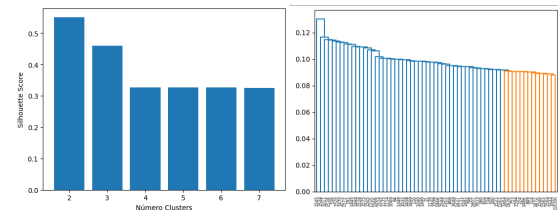


Fig. 6. *Silhouette Scores* e dendrograma (truncado) para *single-linkage*.

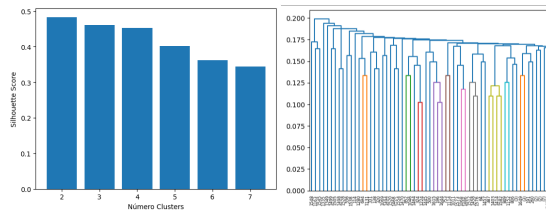


Fig. 7. *Silhouette Scores* e dendrograma (truncado) para *average-linkage*.

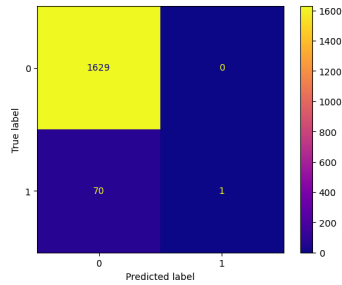


Fig. 8. Matriz de confusão para as *true labels* vs *clusters*.

	Random Index	Adjusted Random Index	Adjusted Mutual Index
Score	0.921	0.026	0.020

TABLE II  
MÉTRICAS PARA DISTÂNCIA DE GOWER

## B. Métodos de Partição

### 1) *K-means*:

## C. DBSCAN

O DBSCAN (*Density-Based Spatial Clustering of Applications with Noise*) é um algoritmo de agrupamento baseado em densidade. Este método descobre *clusters* com forma arbitrária e é bastante eficiente para conjuntos de dados multi dimensionais.

O DBSCAN procura *clusters* através da vizinhança de cada ponto do *dataset* e verifica se contém mais do que um determinado número de observações.

Desta forma, é necessário definir dois parâmetros de entrada *a priori*: o raio de vizinhança (*Eps*) e o número mínimo de pontos para definir um *cluster* (*minPts*).

Pode-se definir três tipos de pontos neste algoritmo. Se uma determinada observação tiver mais pontos do que o número mínimo, definido num raio *Eps*, é considerado um *core point*. Se um ponto tiver menos pontos que número mínimo mas estiver na vizinhança de um *core point* é então definido como um *border point*. Um ponto que não seja nenhum dos anteriores é considerado um *noise point*.

O DBSCAN é capaz de detetar *clusters* de diferentes formas e densidades, e é menos afetado pelo ruído e *outliers* do que outros algoritmos de agrupamento. No entanto, é preciso ser muito rigoroso quando se define valores para o raio e para o número mínimo de pontos, uma vez que uma pequena alteração pode conduzir a resultados muito diferentes.

## IV. CONCLUSÃO

## REFERENCES