



Universidade do Minho

Departamento de Informática

Mestrado em Matemática e Computação

Mestrado em Engenharia de Sistemas

Perfil de Machine Learning: Fundamentos e Aplicações

Sistemas Baseados em Semelhança

1º/2º Ano, 1º Semestre

Ano letivo 2023/2024

Enunciado Prático nº 4

26 de outubro de 2023

Tema

Clustering

Enunciado

Pretende-se, com este enunciado prático, que sejam aplicados métodos de *clustering* sobre um *dataset* de vinhos, o qual contém um ficheiro para aprendizagem e outro para teste. Deverão também ser aplicadas técnicas para exploração e tratamento de dados, assim como para parametrização do *workflow* a desenvolver.

Tarefas

Numa primeira fase devem descarregar o *dataset* disponível em <https://goo.gl/8jjW8t>. Devem, de seguida:

T1. Carregar, no *Knime*, o *dataset* descarregado e explorar os dados;

T2. Tratar os dados, i.e.:

- Fazer cast do atributo “*quality*” para inteiro;
 - Normalizar todos os atributos numéricos utilizando a transformação linear Min-max de forma a produzir um input normalizado entre 0 e 1;
 - Criar 4 *bins* de igual frequência para a *feature* “*citric acid*”, substituindo a *feature* original;
 - Renomear cada *bin* de forma a que o primeiro corresponda a *Low*, o segundo a *Medium*, o terceiro a *High* e o quarto a *Very High*.
- Dica: no passo anterior usar *Numbered* como *Bin Naming* – podem depois usar os nodos *Table Creator* e *Cell Replacer*.

T3. Aplicar:

- Uma Análise de Componentes Principais (PCA) de forma a projetar os dados em apenas duas dimensões;
- Utilizar um *scatter plot* para visualização dos resultados obtidos pelo PCA.

T4. Segmentar o *dataset*:

- Aplicando o método *k-means*;
- Atribuir diferentes cores por qualidade do vinho e diferentes formas aos clusters;
- Criar *scatter plots* e *scatter matrixes* que permitam ter uma noção gráfica, em duas dimensões, dos atributos e dos clusters criados;
- Ler e tratar os dados de teste de forma a que, com base no modelo desenvolvido nos passos anteriores, seja atribuído um cluster a cada registo deste ficheiro;
- Guardar o resultado da atribuição num ficheiro csv.

- T5.** Parametrizar o *workflow*, utilizando variáveis de fluxo para definir o número de *bins*, o número de *clusters* e os títulos dos gráficos criados;
- T6.** Produzir o *workflow* de maneira a que seja possível visualizar, numa única página, todos os componentes visuais implementados;
- T7.** Experimentar, avaliar e comparar outros métodos de segmentação.