



Universidade do Minho

Departamento de Informática

Mestrado [integrado] em Engenharia Informática

Mestrado em Engenharia de Sistemas

Perfil de Machine Learning: Fundamentos e Aplicações

Sistemas Baseados em Similaridade

4º/2º Ano, 1º Semestre

Ano letivo 2020/2021

Enunciado Prático nº 4

05 de novembro de 2020

Tema	<i>Clustering</i>
Enunciado	Pretende-se, com este enunciado prático, que sejam aplicados métodos de <i>clustering</i> sobre um <i>dataset</i> 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 <i>workflow</i> a desenvolver.
Tarefas	<p>Numa primeira fase devem descarregar o <i>dataset</i> disponível em https://goo.gl/8jjW8t. Devem, de seguida:</p> <p>T1. Carregar, no <i>Knime</i>, o <i>dataset</i> descarregado e explorar os dados;</p> <p>T2. Tratar os dados, i.e.:</p> <ol style="list-style-type: none">Fazer cast do atributo “<i>quality</i>” 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 <i>bins</i> de igual frequência para a <i>feature</i> “<i>citric acid</i>”, substituindo a <i>feature</i> original;Renomear cada <i>bin</i> de forma a que o primeiro corresponda a <i>Low</i>, o segundo a <i>Medium</i>, o terceiro a <i>High</i> e o quarto a <i>Very High</i>. <ul style="list-style-type: none">Dica: no passo anterior usar <i>Numbered</i> como <i>Bin Naming</i> – podem depois usar os nodos <i>Table Creator</i> e <i>Cell Replacer</i>. <p>T3. Aplicar:</p> <ol style="list-style-type: none">Uma Análise de Componentes Principais (PCA) de forma a projetar os dados em apenas duas dimensões;Utilizar um <i>scatter plot</i> para visualização dos resultados obtidos pelo PCA. <p>T4. Segmentar o <i>dataset</i>:</p> <ol style="list-style-type: none">Aplicando o método <i>k-means</i>;Atribuir diferentes cores por qualidade do vinho e diferentes formas aos clusters;Criar <i>scatter plots</i> e <i>scatter matrixes</i> 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.