



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 2019/2020

Ficha Prática nº 6

07 de novembro de 2019

Tema	Segmentação/ <i>Clustering</i> .
Enunciado	Criação e análise de <i>clusters</i> sobre dois <i>datasets</i> . O primeiro refere-se a um <i>dataset</i> de vinhos e contém um ficheiro para aprendizagem e um outro para teste (https://goo.gl/8jjW8t). O segundo <i>dataset</i> contém dados referentes a 3 fabricantes de automóveis (https://goo.gl/Eap319).
Tarefas	<p>Para o primeiro <i>dataset</i> o objetivo deste exercício é o de criar clusters de acordo com as características químicas dos vinhos usando algoritmos de segmentação, como o <i>k-means</i>. Deve assim ser desenvolvido um <i>workflow</i> na plataforma <i>Knime</i> para:</p> <ul style="list-style-type: none">• Tratar o atributo <i>quality</i> de forma a torná-lo num inteiro;• Normalizar os atributos numéricos utilizando a transformação linear Min-max de forma a produzir um <i>input</i> normalizado sobre o qual deve ser aplicado o algoritmo de segmentação;• 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 modelo no formato PMML. <p>Para o segundo <i>dataset</i>, deve-se proceder à criação de clusters de forma similar à descrita em cima, i.e.:</p> <ul style="list-style-type: none">• Tratar os atributos e normalizá-los;• Aplicar diferentes algoritmos de segmentação sobre os dados normalizados e sobre os dados não normalizados;• Atribuir diferentes cores e formas aos clusters. Criar <i>plots</i> que permitam ter uma noção gráfica dos atributos. Qual o impacto da normalização dos dados?• Aplicar uma Análise de Componentes Principais (<i>Principal Component Analysis</i>) de forma a projetar os dados em apenas duas dimensões. Qual o impacto da normalização dos dados no PCA?• Guardar o modelo no formato PMML.