



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 3

29 de outubro de 2020

Tema	<i>A Data Science Perspective</i>
Enunciado	Uma exploração aprofundada dos dados permite que se tirem ilações, muitas vezes escondidas, que poderão ser importantes para se compreender o domínio e o problema em mãos. Com este enunciado prático é esperado que sejam aplicadas um conjunto de técnicas que permitam explorar e tratar <i>datasets</i> .
Tarefas	<p>Numa primeira fase devem descarregar o <i>dataset</i> disponível em https://goo.gl/p2y19t que contém dados de um conjunto de utilizadores de uma determinada plataforma web assim como o seu sentimento em relação à mesma. Devem, de seguida:</p> <p>T1. Carregar, no <i>Knime</i>, o <i>dataset</i> descarregado. Aplicar nodos para exploração de dados, i.e., analisar os dados em relação à sua:</p> <ol style="list-style-type: none">Tendência central;Dispersão estatística;Correlação entre <i>features</i>. <p>T2. Criar <i>plots</i> para visualização dos dados;</p> <p>T3. Aplicar nodos para tratamento de dados de forma a:</p> <ol style="list-style-type: none">Excluir todas as colunas do tipo <i>Double</i>;Tratar valores em falta;Remover registos duplicados;Criar 3 <i>bins</i> de igual frequência para a <i>feature age</i>;Para cada registo, extrair o ano, mês e dia da semana da <i>feature birthday</i>;Excluir utilizadores da plataforma que tenham uma atividade na plataforma (<i>WebActivity</i>) inferior a 1 hora e que tenham mais de 70 anos;Excluir todos os registos que contenham a <i>sub-string</i> “co” no produto. <p>T4. Aplicar nodos para agregação de dados de forma a:</p> <ol style="list-style-type: none">Por género, obter o número e a percentagem de registos, assim como a média da idade e da atividade na plataforma. Obter também o mínimo e máximo da idade;Por género e atividade na plataforma, obter a moda da análise do sentimento em relação à plataforma e a média da avaliação do sentimento;

- c. Por análise de sentimento, obter o número de registos, a média do salário anual estimado, o somatório do salário anual e a média do número de contratos.

T5. Análise crítica à informação extraída das agregações efetuadas na tarefa anterior. Que conclusões poderia a empresa tirar?

Numa segunda fase devem descarregar o *dataset* disponível em <https://bit.ly/3525yDr>. Este *dataset* contém dados referentes à performance de vários jogadores de futebol na edição 2017/2018 da *Premier League*. Devem, de seguida:

T6. Carregar, no *Knime*, o *dataset* descarregado. Explorar os dados, procurar informação relevante e mostrar essa mesma informação. P.e., qual a equipa mais indisciplinada? Qual o top-10 dos assistentes para golo? Qual o top-5 de nacionalidades na liga?