



Universidade do Minho

Licenciatura em Engenharia Informática

Mestrado integrado em Engenharia Informática

Aprendizagem e Decisão Inteligentes

3º/4º Ano, 2º Semestre

Ano letivo 2021/2022

Enunciado Prático n.º 7

Abril, 2022

Tema	Seleção de atributos e Validação de Modelos
Enunciado	<p>Uma exploração aprofundada dos dados permite retirar ilações, muitas vezes escondidas, que poderão ser importantes para compreender o domínio e o problema em mãos.</p> <p>Com este enunciado prático os alunos deverão aplicar um conjunto de técnicas que permitam explorar e tratar <i>datasets</i>, por forma a permitir explorar estratégias de seleção de características (<i>feature selection</i>) e usar diversas abordagens de validação de modelos (<i>cross validation</i>, etc.).</p>
Tarefa 1	<p>Numa primeira fase devem descarregar o <i>dataset users_sentiment</i> disponível na plataforma de <i>e-learning</i> da UMinho, secção [Conteúdo]. Este <i>dataset</i> contém dados de um conjunto de utilizadores de uma determinada plataforma <i>web</i> assim como o seu sentimento em relação à mesma.</p> <p>T1. Carregar, no <i>Knime</i>, o <i>dataset users_sentiment</i> e aplicar nodos de exploração de dados como forma de permitir a análise dos dados em relação à sua:</p> <ol style="list-style-type: none">Tendência central;Dispersão estatística;Correlação entre <i>features</i>; <p>tendo sempre em consideração o tipo e respetiva qualidade de dados e quantidade de conhecimento “acondicionadas” no <i>dataset</i>.</p> <p>T2. Obter perspetivas gráficas de análise dos dados criando <i>plots</i> para visualização dos dados.</p> <p>T3. Aplicar nodos para tratamento de dados de forma a, por exemplo:</p> <ol style="list-style-type: none">Excluir todas as colunas do tipo <i>Double</i>;Tratar valores em falta (<i>missing values</i>);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.

Tarefa 2

Na sequência da **Tarefa 1** e utilizando o mesmo *dataset users_sentiment*, aplique um conjunto de estratégias de seleção de características (*feature selection*) e validação de modelos de aprendizagem automática.

Atendendo aos passos aplicados na **Tarefa 1 - T1, T2 e T3**, resolva as seguintes tarefas:

T1. Aplicar nodos de tratamento de dados de forma a permitir o treino e validação de modelos de aprendizagem, incluindo, mas não se limitando a:

- a. Transformar os valores categóricos em representações de valores numéricos;
- b. Remover da análise de dados colunas potencialmente indesejadas;
- c. Aplicar estratégias de normalização de dados a colunas numéricas que entenda relevante;
- d. Tratar casos que sejam considerados *outliers*.

T2. Construir modelos de decisão, como, por exemplo:

- a. Árvore de decisão para prever o valor de *Sentiment Analysis*;
- b. Aplicar técnicas de validação de modelos (*k-fold cross-validation*), obtendo a tabela com os resultados de validação do modelo para cada *fold*;
- c. Explorar o treino e validação de outras técnicas de modelos de aprendizagem supervisionadas.

T3. Usar os resultados obtidos, em particular, com as técnicas de *feature selection*, para desenvolver novos modelos de aprendizagem e comparar os diferentes modelos.