



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 4

Março, 2022

Tema	Regressão Linear
Enunciado	Regressão linear é uma técnica de <i>machine learning</i> usada para estimar o objetivo de um caso de estudo, sendo dado um conjunto de dados que caracterizam o comportamento de uma série de casos passados.
Tarefas	<p>Numa primeira fase devem descarregar o <i>dataset</i> [lego_sets] disponível na plataforma de <i>e-learning</i> da UMinho, secção [Conteúdo]. LEGO é uma marca de blocos de construção e são, normalmente, vendidos em conjuntos para construir um objeto específico. Este <i>dataset</i> contém dados sobre esses conjuntos, incluindo o número de peças de cada conjunto, preço de venda, idade sugerida, entre outros. O objetivo deste exercício foca-se no desenvolvimento de um modelo de regressão linear capaz de prever o preço do <i>set</i> (i.e., 'list_price'), dado um conjunto de <i>features</i> disponibilizados no <i>dataset</i>.</p> <p>T1. Carregar, no <i>Knime</i>, o <i>dataset</i> [lego_sets] e aplicar nodos de exploração de dados como forma de permitir a análise dos dados;</p> <p>T2. Proceder ao tratamento e limpeza dos dados:</p> <ol style="list-style-type: none">Excluir o conjunto de <i>features</i> que considere irrelevantes para o desenvolvimento e validação do modelo de aprendizagem;Aplicar técnicas de <i>feature engineering</i> e <i>data encoding</i> nas respetivas <i>features</i> do <i>dataset</i>, incluindo:<ol style="list-style-type: none">Técnicas de <i>encoding</i>;Tratamento de valores em falta;Tratamento de <i>outliers</i>;Tratamento de casos repetidos. <p>T3. Analise a correlação das <i>features</i> através de uma matriz de correlação. Dada esta informação, que <i>features</i> sugeria remover?</p> <p>T4. Particionar os dados utilizando 80% para treino e 20% para teste;</p> <p>T5. Aplicar um nodo <i>Linear Regression Learner</i> para treinar um modelo de regressão e um nodo <i>Linear Regression Predictor</i> para obter previsões utilizando o modelo treinado;</p> <p>T6. Avaliar o desempenho da previsão do modelo de regressão utilizando o nodo <i>Numeric Scorer</i>;</p> <p>T7. Criar um <i>scatter plot</i> como forma de observar os valores previstos <i>versus</i> valores reais;</p> <p>T8. Calcular o erro residual e analisar a sua distribuição;</p> <p>T9. Atendendo aos resultados e conhecimentos adquiridos ao longo da ficha prática, explore a utilização de outras técnicas de tratamento e limpeza dos dados, de modo a minimizar o erro do modelo de aprendizagem.</p>