



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 5

Março, 2022

Tema	Modelos de Aprendizagem baseados em Árvores de Decisão
Enunciado	Com este enunciado prático pretende-se que sejam desenvolvidos modelos de aprendizagem baseados em árvores de decisão, abordando parâmetros nominais e numéricos como a medida de qualidade, o método de <i>pruning</i> e o número mínimo de registos por nodo, entre outros.
Tarefas	<p>Uma multinacional na área do retalho possui o histórico de vendas semanais de 17 lojas em diferentes regiões do país, sendo que cada loja contém vários departamentos (desporto, cozinha, produtos alimentícios, higiene pessoal, entre outros). A empresa realiza também vários eventos promocionais ao longo do ano, normalmente precedendo feriados importantes.</p> <p>A empresa pretende agora extrair informação relevante dos <i>datasets</i> e desenvolver um modelo de aprendizagem que, com base num conjunto relevante de <i>features</i>, permita estimar as vendas mensais de cada loja. A empresa possui dois <i>datasets</i>: o primeiro (<i>dataset</i> [lojas] disponível na plataforma <i>e-learning</i> da UMinho, secção [Conteúdo]) contém informação sobre cada uma das lojas, incluindo o seu tipo e tamanho; o segundo (<i>dataset</i> [sales_training_set] disponibilizado de modo idêntico) contém dados referentes às vendas semanais de cada departamento de cada loja, a data e um <i>boolean</i> indicando se houve um feriado durante essa semana.</p> <p>Um terceiro dataset (<i>dataset</i> [sales_test_set] disponível nos mesmos termos) deve ser utilizado, única e exclusivamente, como conjunto de teste por ocasião do desenvolvimento dos modelos de aprendizagem, de modo a garantir que estes são avaliados com dados que desconhecem.</p> <p>Deve ser desenvolvido um <i>workflow</i> para:</p> <p>T1. Carregar no <i>Knime</i> os dois primeiros <i>datasets</i>, juntá-los e explorar os dados utilizando nodos de visualização gráfica que permitam interpretar a análise efetuada;</p> <p>T2. Proceder ao tratamento e limpeza dos dados:</p> <ol style="list-style-type: none">Fazer <i>label encoding</i> da <i>feature isHoliday</i> (o valor 1 deverá corresponder ao valor <i>True</i>);Adicionar, a cada registo, as <i>features</i> ano e mês;Agrupar os registos por loja, tipo, tamanho, ano e mês, agregando de forma a obter o somatório de vendas semanais de cada loja e a indicação da existência de feriados nesse mês;Normalizar o somatório das vendas semanais utilizando a transformação linear <i>Min-Max</i> entre 0 e 1;Criar 4 <i>bins</i> de igual frequência sobre o valor normalizado no passo anterior (ativar a opção <i>replace target column(s)</i>);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>;

T3. Treinar:

- a) Uma árvore de decisão capaz de prever o valor de vendas de cada mês para cada uma das 17 lojas;
- b) Carregar o *dataset* de teste e prever o valor de vendas de cada mês para cada uma das 17 lojas;
- c) Mostrar, graficamente, uma tabela com a matriz de confusão do modelo;

T4. Fazer o *tuning* do modelo criado no passo anterior, experimentando:

- a) Todos os valores, entre 2 e 10, para o número mínimo de registros por nodo;
- b) Todas as possibilidades para a medida de qualidade;
- c) Todas as possibilidades para o método de *pruning*;
- d) Guardar e analisar todos os resultados obtidos para cada combinação de hiper-parâmetros averiguada. Qual a combinação que oferece melhor desempenho? Existem grandes discrepâncias?

T5. Treinar um modelo de aprendizagem de floresta aleatória (*random forest*). Guardar e analisar todos os resultados obtidos para cada combinação de hiperparâmetros averiguada;

T6. Analisar e comparar os desempenhos dos modelos treinados anteriormente. Que conclusões se podem tirar?