



Universidade do Minho

Departamento de Informática

Mestrado [integrado] em Engenharia Informática

Dados e Aprendizagem Automática

1º/4º Ano, 1º Semestre

Ano letivo 2023/2024

Enunciado Prático nº 6

2 de novembro de 2023

Tema	Máquinas de Vetores de Suporte e Árvores de Decisão: Ajuste de Hiperparâmetros com <i>GridSearchCV</i> e <i>Pruning</i>
Enunciado	<p>Máquinas de vetores de suporte são modelos lineares de aprendizagem supervisionada aplicados a problemas de classificação ou de regressão, onde a técnica cria hiperplanos que separam os casos de estudo em classes. O objetivo é encontrar um hiperplano que tenha a margem máxima, ou seja, a distância máxima entre os pontos de dados de ambas as classes. Este modelo permite resolver problemas lineares e não lineares, produzindo modelos com <i>accuracy</i> significativa e com menor esforço computacional.</p> <p>Árvores de decisão são modelos lineares de aprendizagem supervisionada também aplicados a problemas de classificação ou de regressão. Nesta técnica, os dados são continuamente divididos de acordo com um determinado parâmetro, onde as “folhas” representam as decisões ou resultados estimados.</p> <p><i>Grid search</i> é uma técnica de ajuste de hiperparâmetros que pode facilitar a construção de um modelo e a avaliação de um modelo para cada combinação de parâmetros de algoritmos por grelha.</p> <p><i>Pruning</i> é uma técnica de remoção de folhas e subárvores de uma árvore de decisão aplicada quando estas partes não apresentam significativa contribuição para a precisão e/ou interpretabilidade da árvore, reduzindo-se assim a complexidade da árvore e aumenta-se a sua generalização.</p>
Tarefas	<p>Neste enunciado, utilizaremos o <i>dataset</i> de disponibilizado (<i>incidents.csv</i>), que contém informações de multi-tipo. Para o desenvolvimento de um modelo de classificação, foi decidido aplicar-se os modelos máquina de vetores de suporte e árvores de decisão.</p> <p>Atendendo ao problema em questão, deverão seguir os seguintes passos:</p> <p>T1. Carregar o <i>dataset</i>, utilizando a função <i>pandas.read_csv(...)</i>;</p> <p>T2. Aplicar métodos para exploração e visualização de dados;</p> <p>T3. Preparar e organizar os conjuntos de casos de estudo do <i>dataset</i> em dados de treino e teste, utilizando a função <i>sklearn.model_selection.train_test_split(..., test_size = 0.3)</i>;</p> <p>T4. Treinar um modelo de máquina de vetores de suporte (<i>sklearn.svm.SVC</i>) e um modelo de árvore de decisão como classificador (<i>sklearn.tree.DecisionTreeClassifier</i>).</p> <p><i>Nota:</i> Definir o X e o y. Atenção ao tipo dos atributos que fazem parte do X;</p> <p>T5. Obter matrizes de confusão e relatório de classificação dos modelos e efetuar a respetiva análise crítica. Avaliar a <i>accuracy</i> do modelo na previsão de ‘<i>incidents</i>’. Avaliar também o modelo usando a métrica <i>f1_macro</i>. Que conclusões se podem tirar?</p> <p>T6. Aplicar a técnica de <i>gridsearchCV</i> (<i>sklearn.model_selection.GridSearchCV</i>) como forma de procurar o conjunto de hiperparâmetros capaz de otimizar o desempenho da classificação dos modelos de máquina de vetores de suporte (<i>C</i> e <i>gamma</i>) e de árvore de decisão (<i>criterion</i> e <i>max_depth</i>). Qual a variação no desempenho do modelo subjacente a estas alterações?</p>

T7. Aplicar a técnica de *pruning* como forma de identificar o conjunto de hiperparâmetros capaz de otimizar o desempenho da classificação do modelo de árvore de decisão (*max_depth* e *ccp_alphas*). Qual a variação no desempenho do modelo subjacente a estas alterações?