

Algoritmos de Árvore de Decisão e Random Forest

Olá! Seja bem-vindo(a) à nossa newsletter semanal sobre Ciência de Dados!

Nesta edição, vamos explorar dois algoritmos fundamentais em ciência de dados: Árvores de Decisão e Random Forest. Esses métodos são amplamente utilizados em tarefas de classificação e regressão devido à sua simplicidade e eficácia.

O que é uma Árvore de Decisão? Uma árvore de decisão é um modelo de aprendizado de máquina que utiliza uma estrutura de árvore para tomar decisões com base nas características dos dados. Cada nó da árvore representa uma decisão baseada em um atributo, e cada ramo representa o resultado dessa decisão, levando a outros nós ou a uma folha que representa a classificação final ou valor previsto. Seus princípios de funcionamento são: 1) **Divisão Recursiva:** A árvore é construída dividindo-se iterativamente os dados em subconjuntos com base em testes de atributos; 2) **Critério de Divisão:** Utiliza critérios como Gini, entropia ou variância para escolher o melhor atributo para a divisão em cada nó; e 3) **Podar:** Para evitar o *overfitting*, a árvore pode ser podada removendo ramos que fornecem pouca informação. Uma árvore de decisão é classificada como um algoritmo de aprendizado supervisionado, usado tanto para classificação quanto para regressão.

O que é Random Forest? Random Forest é um algoritmo ensemble que constrói múltiplas árvores de decisão durante o treinamento e usa a média ou a votação majoritária dos resultados dessas árvores para melhorar a precisão e controlar o *overfitting*. Seus princípios de funcionamento são: 1) **Bootstrap Aggregating (Bagging):** Cria várias subamostras dos dados de treinamento com reposição e constrói uma árvore de decisão para cada subamostra; 2) **Randomização de Atributos:** Cada árvore é construída a partir de um subconjunto aleatório de atributos, aumentando a diversidade entre as árvores; e 3) **Agregação:** Os resultados das árvores são combinados por votação (para classificação) ou média (para regressão). Random Forest é classificado como sendo um algoritmo ensemble, que combina múltiplas árvores de decisão.

A seguir, uma tabela compara esses dois métodos/algoritmos:

Método/Algoritmo	Resumo do Método	Onde Utilizar	Onde Não Utilizar	Exemplo de Código Python
Árvores de Decisão	Modelo baseado em estrutura de árvore para tomar decisões e prever valores com base em atributos dos dados.	Tarefas de classificação e regressão com dados estruturados e interpretáveis.	Dados com muitos ruídos, onde pode ocorrer <i>overfitting</i> , e problemas que exigem alta precisão.	<pre>from sklearn.tree import DecisionTreeClassifier clf = DecisionTreeClassifier() clf.fit(X_train, y_train) y_pred = clf.predict(X_test)</pre>
Random Forest	Algoritmo ensemble que combina múltiplas árvores de decisão para melhorar a precisão e evitar <i>overfitting</i> .	Quando se deseja alta precisão e controle sobre o <i>overfitting</i> , em tarefas de classificação e regressão.	Em cenários onde a interpretabilidade do modelo é crucial, devido à complexidade do ensemble.	<pre>from sklearn.ensemble import RandomForestClassifier clf = RandomForestClassifier(n_estimators=100) clf.fit(X_train, y_train) y_pred = clf.predict(X_test)</pre>

Os algoritmos de árvore de decisão e Random Forest são ferramentas poderosas no arsenal de um cientista de dados. Enquanto as árvores de decisão oferecem simplicidade e interpretabilidade, o Random Forest fornece robustez e precisão através de uma abordagem ensemble. Entender as nuances de quando e como utilizar esses métodos é crucial para maximizar seu impacto em projetos de análise de dados.

Esperamos que você tenha achado estas informações úteis. Seja você um iniciante no mundo da Ciência de Dados ou um profissional experiente, dominar esses algoritmos pode aumentar suas chances de sucesso em projetos e na carreira. Se você tiver alguma dúvida ou sugestão, não hesite em entrar em contato. Queremos garantir que nossos conteúdos sejam sempre relevantes e úteis para você.

Fique atento à nossa próxima edição, onde continuaremos a explorar ferramentas técnicas e dicas valiosas para aprimorar suas habilidades em Ciência de Dados. Até lá, continue explorando e aprendendo!

Saudações,

Prof. Dr. Dilermando Piva Jr

Coordenador de Ciência de Dados para Negócios / Fatec Votorantim

E-mail: f301.cdn@fatec.sp.gov.br