

Data 610

Chapter 9

9.1 Introduction to CART

Classification and Regression Trees are a set of decision-making tools that organize data into a tree-like structure of decisions and outcomes. Imagine you're trying to predict whether to bring an umbrella when you leave the house. You might consider factors like the weather forecast (sunny, cloudy, rainy), the day of the week (weekday or weekend), or the season. CART helps in systematically breaking down these decisions to predict an outcome (bring an umbrella or not) based on historical data. This method is intuitive and mimics human decision-making processes, making it a powerful tool for predictions.

9.2 Classification Trees Detailed

- **Recursive Partitioning:** This is a divide-and-conquer approach where the data is split into increasingly smaller subsets to create a decision tree. Think of it as organizing a large group of people into smaller, more manageable groups based on shared characteristics until each group is as homogeneous as possible.
- **Categorical Predictors:** These predictors are variables that can take on one of a limited set of classes. For instance, if you're predicting the success of a marketing campaign (success or failure), categorical predictors might include the type of campaign (email, social media, billboard).
- **Measures of Impurity:** Imagine you have a bag of colored balls, and you want to organize them so that each bag contains balls of only one color. Measures like Gini impurity or entropy help you decide the best way to divide the balls

to achieve this goal. In the context of classification trees, these measures help determine the best way to split the data to make the subsets as pure as possible (i.e., containing instances of a single class).

- **Tree Structure:** The structure consists of nodes (questions or decisions), branches (possible answers), and leaves (final outcomes or predictions). For example, a node could represent a question about the age of a customer, branches could represent age groups, and leaves could represent the likelihood of a customer making a purchase.
- **Classifying a New Record:** This involves navigating through the tree based on the attributes of the new record until reaching a leaf, which provides the classification. It's like following a recipe where each step (node) asks a question about the ingredients you have, and depending on your answer, you're guided to the next step until you reach the final dish (classification).

9.3 Evaluating the Performance of Classification Trees

The tree's predictive performance is assessed using metrics such as accuracy (the proportion of correct predictions out of all predictions), precision (the proportion of true positive predictions in all positive predictions), and recall (the proportion of true positive predictions out of all actual positives). These metrics help determine how well the tree generalizes its predictions to new, unseen data.

9.4 Avoiding Overfitting in Depth

- **Stopping Tree Growth:** Conditional inference trees are a sophisticated method that uses statistical tests to decide whether to continue splitting a node. It's a way of asking, "Is this split statistically significant?" If not, the tree stops growing to prevent overfitting.

- **Pruning The Tree:** After a tree has grown fully, it might contain branches that make the model too complex. Pruning removes these branches, simplifying the model without significantly reducing accuracy.
- **Cross-Validation:** This involves dividing the data into several parts, using some for training the tree and the rest for testing it. This process is repeated multiple times to ensure that the model performs well across different subsets of the data.
- **Best-Pruned Tree:** The ideal version of the tree is one that achieves the best balance between accuracy and simplicity. This tree is identified through cross-validation, ensuring it performs well on unseen data without being overly complex.

9.5 Classification Rules from Trees

Rules extracted from trees are straightforward if-then statements that describe the path from the root to a leaf. For example, "If the customer is under 30 and has made a purchase in the last month, then they are likely to respond to the new marketing campaign." These rules are valuable for making decisions transparent and understandable.

9.6 Classification Trees for More Than Two Classes

In scenarios with more than two classes, classification trees can still be applied. For example, if predicting the type of pet someone will choose (dog, cat, bird), the tree can handle this multi-class problem by finding the best splits to segregate the data into these distinct categories.

9.7 Regression Trees Explained Further

- **Prediction:** Regression trees predict a continuous value. For example, predicting a house's sale price based on features like square footage, number of bedrooms, and age of the house.

- **Measuring Impurity:** In regression trees, impurity is measured by the variance or mean squared error within each node. The goal is to find splits that result in subsets of data with as little variance as possible in their outcomes.
- **Evaluating Performance:** The performance of regression trees is typically evaluated using metrics like Mean Squared Error (MSE) or Root Mean Squared Error (RMSE), which measure the average difference between the predicted and actual values.

9.8 Improving Prediction with Ensemble Methods

- **Random Forests:** By creating a 'forest' of many decision trees from random subsets of the data and averaging their predictions, random forests improve prediction accuracy and robustness. It's like asking a crowd of experts for their opinion rather than relying on just one.
- **Boosted Trees:** Boosting involves sequentially building trees, where each tree tries to correct the errors of the previous one. This process continues until the model's performance can no longer improve. It's akin to a team working together, with each member addressing and fixing the mistakes of the previous one to achieve the best result.

9.9 Advantages and Weaknesses of Trees in Detail

Trees are highly interpretable and can easily handle qualitative (categorical) and quantitative (numerical) data. They require relatively little data preparation. However, they can be prone to overfitting and might not capture complex relationships or interactions between variables as well as some other models, like neural networks.

Conclusion and Implications for Learning

Understanding CART is crucial for anyone entering the field of data science because these models are foundational to many types of predictive analysis across various domains. They offer a

balance between simplicity and predictive power, making them an excellent starting point for tackling complex real-world problems. By learning these concepts, you equip yourself with the tools to make informed decisions based on data, a skill highly valued across industries.