

Analyzing a medical dataset containing patient information and drug response patterns provides valuable insights into the effectiveness of decision tree algorithms in healthcare predictions. The dataset, obtained from Kaggle(via IBM), comprises 200 patient records with various health metrics including age, blood pressure, cholesterol levels, and sodium-to-potassium ratios, making it an ideal candidate for exploring classification techniques. The choice of medical data is particularly relevant as it represents real-world healthcare scenarios where accurate prediction of drug responses can significantly impact patient outcomes.

Working with this dataset presented unique opportunities to explore different decision tree implementations. The data's clean structure, with no missing values and a mix of categorical and numerical variables, provided an excellent foundation for testing various tree-based approaches. The features included both basic patient demographics and specific health indicators, allowing for comprehensive analysis of factors influencing drug response patterns.

The analysis employed three distinct approaches: two different decision tree configurations and a random forest model. Initial data preparation involved converting categorical variables to factors and implementing a strategic 80-20 split for training and testing, ensuring proportional representation of drug classes. This preparation was crucial for maintaining the integrity of our classification models and ensuring reliable predictions.

Interestingly, all three models achieved identical performance metrics, with an accuracy of 97.37% and a Kappa value of 0.9611. The consistency in performance across different model configurations revealed strong, stable patterns in the data. The first decision tree identified the sodium-to-potassium ratio as the primary splitting criterion, followed by blood pressure and age as secondary decision nodes. Attempts to force alternative splitting arrangements in the second tree resulted in the same structure, suggesting the robustness of these relationships.

The random forest model, despite its more complex ensemble approach, maintained the same level of accuracy as the individual decision trees. This consistency across models indicates that the relationships between patient characteristics and drug responses are well-defined and can be effectively captured by simpler tree structures. The variable importance plot from the random forest confirmed the sodium-to-potassium ratio as the strongest predictor, while patient sex emerged as the least influential feature.

From a medical perspective, the single decision tree model would be the preferred choice for implementation. It offers equivalent performance to the more complex random forest while providing clear, interpretable decision paths that healthcare providers can easily follow. The tree's structure reveals logical splitting points based on clinical measurements, making it particularly valuable for practical medical decision-making.

The analysis highlights the importance of feature selection in medical predictions, as demonstrated by the consistent significance of the sodium-to-potassium ratio across all models. However, it also reveals that additional model complexity doesn't necessarily yield better results when the underlying relationships are strong and well-defined. The perfect prediction rates for three of the five drug classes, with only a single misclassification between two specific drugs, suggests that the model has captured meaningful patterns in patient response to different medications.

This study demonstrates that decision trees can effectively predict drug responses while maintaining interpretability, a crucial factor in medical applications. The consistency in performance across different tree-based approaches underscores the stability of the identified patterns, providing confidence in the model's potential real-world application for drug prescription guidance.