

# **the implementation of a machine learning algorithm based on a decision tree model**

In this report, we present the implementation of a machine learning algorithm based on a **decision tree model**. This algorithm is designed to analyze and predict class labels in a given dataset, specifically focusing on the classification task. The dataset used in this implementation is loaded from an Excel file named **injury2.xlsx**. The goal is to create a robust decision tree model that can accurately classify new instances based on the learned patterns from the training data. Below, we provide a detailed step-by-step explanation of the code and its functionality.

## **Data Loading and Preparation**

### **1. Loading Data:**

- The dataset is read from the injury2.xlsx file using the xlsread function, and the data is stored in three variables: num (numeric data), str (string data), and row (all data).
- The first row, which contains column headers, is removed from row, and the corresponding labels are extracted and stored in the label variable.

### **2. Data Conversion:**

- The code checks for numeric values stored as strings in the tenth column. If such values are found, they are converted to numeric format and updated in both num and row.

## **Feature Identification**

### **3. Identifying Numeric and Categorical Features:**

- The code differentiates between numeric and categorical features by examining the first row of the dataset.
- Numeric features are identified and stored in the o1 array, while categorical features are identified and stored in the o2 array.

## **Data Splitting**

### **4. Splitting Data into Training and Testing Sets:**

- The dataset is divided into five parts to facilitate cross-validation. Each part consists of 20% testing data and 80% training data.
- This is achieved by randomly selecting rows for the testing set and ensuring that the same rows are not repeatedly selected across different parts.

## **Identifying Unique States**

### **5. Identifying Unique States of Categorical Features:**

- For each categorical feature, the code identifies the unique states or categories present in the data.
- These unique states are counted and stored in the `halat` array for further processing.

## Class Separation

### 6. Dividing Training Data into Two Classes:

- The training data is divided into two subsets based on the class labels extracted earlier.
- This separation is necessary for building class-specific models.

## Building Probability Tables

### 7. Building Probability Tables:

- For each categorical feature, probability tables are constructed to store the occurrence counts of different states for each class.
- These tables are used to calculate the likelihood of each state occurring within a particular class.

## Decision Tree Construction

### 8. Constructing the Decision Tree Model:

- The decision tree model is built for each part of the data. The model includes separate trees for both categorical and numeric features.
- For categorical features, a tree structure is created based on the probability tables.
- For numeric features, thresholds are determined to split the data, and the tree is constructed accordingly.

### 9. Evaluating the Model:

- The model's performance is evaluated using the training data. This evaluation involves calculating the accuracy and error rates to ensure the model is learning correctly.

## Prediction

### 10. Predicting Labels:

- The constructed decision tree model is used to predict the labels for the testing data.
- The predictions are compared with the actual labels to assess the model's accuracy.

## Error Calculation

### 11. Calculating Error:

- The prediction error for each part is calculated, and the mean error across all parts is computed.
- This mean error provides an overall measure of the model's performance.

## Summary and Variables

Finally, several key variables are saved for further analysis:

- `mianginkhata`: Mean prediction error across all parts.
- `finaltree`: The final decision tree model.
- `predictlabel`: Predicted labels for the testing data.
- `actuallabel`: Actual labels for the testing data.
- `testeror`: Error rates for each part of the data.

These variables are crucial for evaluating the model's performance and can be used for further refinement and analysis.

## Conclusion

In summary, the implemented code provides a comprehensive solution for building and evaluating a decision tree model for classification tasks. By carefully preparing the data, identifying feature types, and using cross-validation, the model aims to achieve high accuracy and reliability in predicting class labels for new instances. This detailed approach ensures that the model is robust and capable of handling both numeric and categorical features effectively.