# Mathematical Essay on Random Forest

Sparsh Tewatia
*Department of Data Science, Joint MS(UoB)*
*I.I.T Madras*
Chennai, India
sparshteotia5@gmail.com

*Abstract*—This assignment presents a mathematical essay on the Random Forest classifier, structured according to the IEEE conference style format. The primary objectives of this work are threefold: first, to cultivate expertise in effectively communicating technical topics; second, to elucidate the mathematical concepts underpinning Random Forest; and third, to employ Random Forest to address a real-world problem.

*Index Terms*—Random Forest, Car-Safety Evaluation, Machine Learning, Analysis

## I. Introduction

In this section, we provide a comprehensive introduction to the assignment, serving as a foundation for the subsequent exploration of the Random Forest classifier and its application to a real-world problem. The introduction is organized into four paragraphs, each addressing a specific aspect of our assignment.

### A. Overview of the Assignment

In this assignment, we aim to develop expertise in writing and communicating about technical topics. This will be achieved by adhering to the IEEE conference style format throughout the paper. Furthermore, we endeavor to elucidate the mathematical ideas that underlie the Random Forest classifier and apply this knowledge to understanding a real-world problem. The assignment takes the form of a mathematical essay on the Random Forest classifier.

### B. Technical Aspects: Random Forest

The Random Forest classifier, a versatile machine learning algorithm, is at the heart of this assignment. This subsection will provide an overview of the key principles and mathematical foundations of the Random Forest algorithm. Understanding these technical aspects is crucial for its application in addressing real-world challenges.

### C. The Problem: Car Classification based on Safety

The problem we seek to address revolves around the classification of cars based on safety attributes. We will use the Car Evaluation Database, which originated from a hierarchical decision model. Accurate classification of cars concerning safety is vital for consumers, manufacturers, and policymakers. In this subsection, we present the problem in detail, including the dataset description and its significance.

### D. Overview of the Essay

This essay is organized into several sections to provide a comprehensive understanding of random forest and their application in solving the car safety classification problem. We will start by delving into the fundamental principles of random forest in Section 2. In Section 3, we will introduce the Car Evaluation data-set, which serves as the basis for our analysis. Section 4 will outline the problem and visualize the data, demonstrating the application of random forest techniques. Finally, in Section 5, we will summarize our key findings and propose potential directions for further research.

Overall, our study demonstrates the potential of random forest as a tool for classifying safety using car evaluation data.

## II. Random Forest

Random Forest is a powerful ensemble learning method widely used in machine learning for classification and regression tasks. It combines the predictions of multiple decision trees to enhance predictive accuracy and reduce overfitting.

### A. Key Principles of Random Forest

*1) Decision Trees:* At the core of Random Forest are decision trees. A decision tree is a hierarchical structure that recursively partitions the data into subsets based on the values of input features. The decision at each node is made by minimizing impurity, often measured by Gini impurity or information gain. A decision tree can be represented as a function $f(x)$ that maps input features $x$ to a predicted outcome.

*2) Bootstrapping:* Random Forest leverages bootstrapping to create multiple subsets of the training data. Bootstrapping involves random sampling with replacement from the training data, resulting in datasets of the same size as the original but with some data points repeated and others omitted. Each of these bootstrapped datasets is used to train a separate decision tree.

*3) Random Feature Selection:* Another critical element is random feature selection. At each split in a decision tree, Random Forest only considers a random subset of features. This randomness reduces the correlation between individual trees and enhances the diversity of the ensemble. The idea is to ensure that no single feature dominates the decision-making process.

*4) Voting or Averaging:* For classification tasks, Random Forest combines the predictions of individual decision trees by majority voting. Each decision tree casts a vote for the class it predicts, and the class with the most votes becomes the final prediction. In regression tasks, the output of each tree is averaged to produce the final prediction.

## B. Mathematical Foundations

The mathematical foundation of Random Forest can be represented as an ensemble of decision trees. Suppose we have a Random Forest model with $N$ decision trees. The overall prediction of the Random Forest for a given input $x$ is the average (for regression) or majority vote (for classification) of the individual tree predictions:

$$F(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x)$$

Where:

$F(x)$ - The Random Forest model's prediction for input $x$

$N$ - The number of decision trees in the ensemble

$f_i(x)$ - The prediction of the $i$th decision tree for input $x$

Each decision tree, $f_i(x)$, makes predictions based on a subset of the training data and a random subset of features. The final prediction is an aggregation of the individual tree predictions, which allows Random Forest to approximate complex functions by combining simpler decision trees.

## C. Advantages of Random Forest

Random Forest offers several advantages, making it a popular choice in various machine learning applications:

- Improved accuracy: The ensemble nature of Random Forest typically results in better predictive accuracy compared to individual decision trees.
- Reduced overfitting: By combining multiple trees, overfitting is mitigated, making Random Forest robust to noisy or complex data.
- Robustness: It can handle high-dimensional data and variables with different scales.
- Feature importance: Random Forest provides a feature importance ranking, helping with interpretability and model understanding.

In Section 4, we will apply the principles of Random Forest to address the specific problem of car classification based on safety using the Car Evaluation Database.

## III. DATA

### A. Dataset Description

The dataset used in this study is the "Car Evaluation Database," derived from a simple hierarchical decision model. This dataset serves as the foundational data source for our analysis and decision tree-based classification. It comprises several key features and a target variable for car safety assessment.

The dataset features are as follows:
- **Buying Price:** The buying price of the car, categorized as "vhigh" (very high), "high," "med" (medium), or "low."
- **Maintenance Price:** The price of car maintenance, categorized as "vhigh," "high," "med," or "low."
- **Number of Doors:** The number of doors in the car, categorized as 2, 3, 4, 5, or "more."
- **Capacity in Terms of Persons to Carry:** The car's passenger capacity, categorized as 2, 4, or "more."
- **Luggage Boot Size:** The size of the luggage boot, categorized as "small," "med" (medium), or "big."
- **Estimated Safety of the Car:** The estimated safety level of the car, categorized as "low," "med" (medium), or "high."

### B. Target Variable

The target variable in our analysis is "Safety." It represents the safety level classification of cars and can take on values such as "unacc" (unacceptable), "acc" (acceptable), "good," and "vgood" (very good).

This dataset will be used to train and test our decision tree-based classification model, with the goal of predicting the safety level of cars based on their attributes.

In the following section, we will delve into the specific problem we aim to address using this dataset and decision tree techniques.

## IV. THE PROBLEM AND ANALYSIS

In this section, we delve into the problem of classifying cars based on their safety features using a random forest. We explore the dataset, conduct data preprocessing, and apply decision tree techniques to gain insights into the decision-making process for car safety assessment.

### A. Problem Statement

The central problem we aim to address in this study is the classification of cars into different safety categories. The safety of a car is a critical factor for consumers, manufacturers, and regulatory authorities. Therefore, developing an effective and interpretable model for car safety assessment is of utmost importance.

We formulate the problem as follows: Given a set of car attributes, including buying price, maintenance price, number of doors, passenger capacity, luggage boot size, and estimated safety, our goal is to predict the safety category of the car. The safety categories are categorized as "unacc" (unacceptable), "acc" (acceptable), "good," and "vgood" (very good).

### B. Data Preprocessing

To prepare the dataset for modeling, we performed several data preprocessing steps:

- We loaded the Car Evaluation Database and explored its structure and features. - Categorical features were one-hot encoded to transform them into a numeric format suitable for machine learning. - The dataset was split into training and testing sets, ensuring the model's ability to generalize to unseen data.

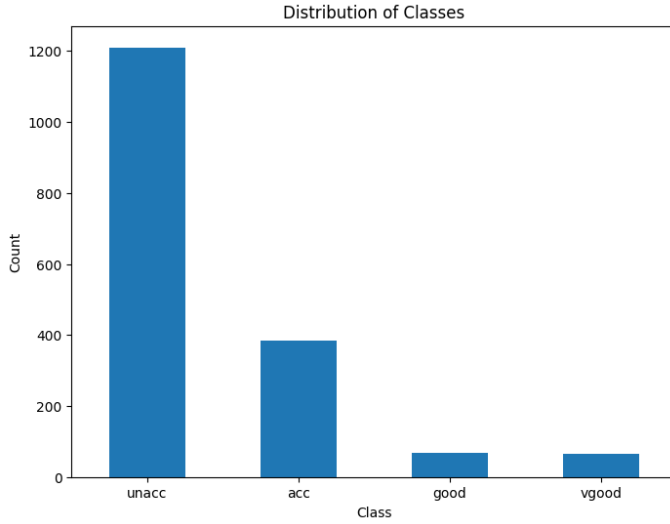We get following class distribution for target variable. (Fig 1)



Fig. 1. Count Plot for class

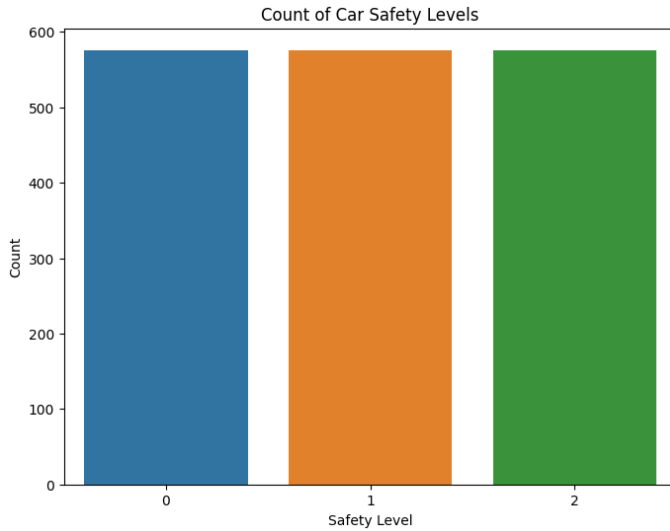We also check the count of car safety level to get the sense of data



Fig. 2. Count of safety level

## C. Random Forest Application

We applied the decision random forest classifier to the preprocessed data for car safety classification. The decision tree algorithm learned from the training data.

## D. Insights and Observations

The application of the random forest algorithm to the car safety classification problem has provided us with valuable insights into the decision-making process.

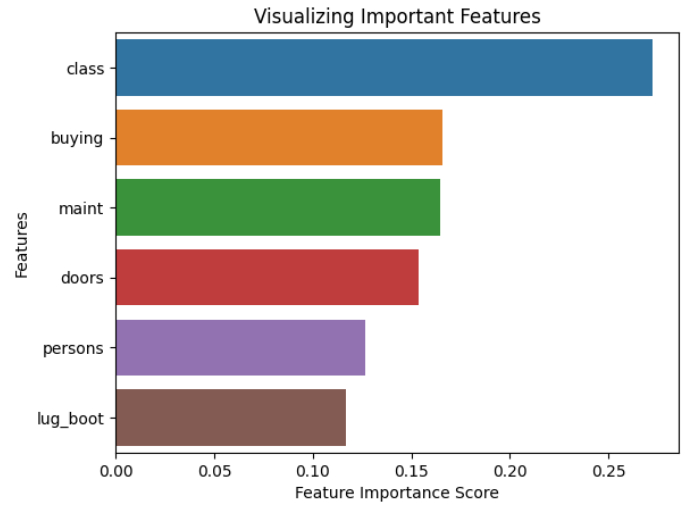We get the feature importance on test data as shown in Fig 3.



Fig. 3. Feature importance

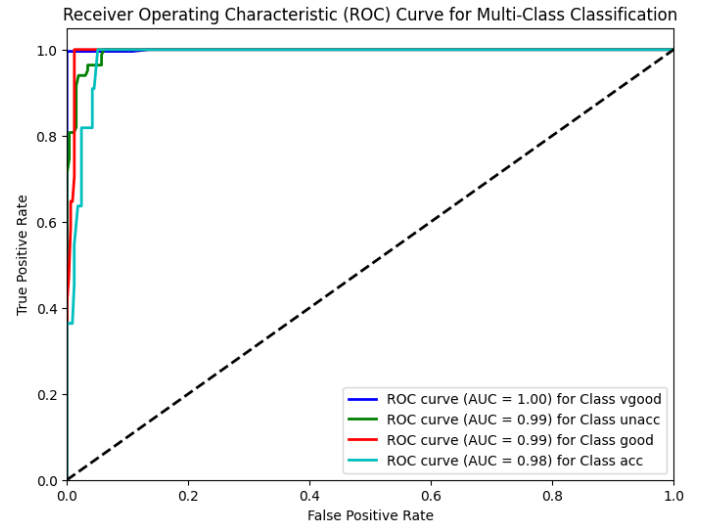We get following ROC curve for the data for each class as shown in Fig 4.



Fig. 4. ROC Curve

## E. Classification Report

We get following classification report.

| Class | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| 0 | 0.99 | 1.00 | 0.99 | 235 |
| 1 | 0.96 | 0.81 | 0.88 | 83 |
| 2 | 0.75 | 0.71 | 0.73 | 17 |
| 3 | 0.39 | 0.82 | 0.53 | 11 |
| **Accuracy** | | | 0.93 | 346 |
| **Macro Avg** | 0.77 | 0.83 | 0.78 | 346 |
| **Weighted Avg** | 0.95 | 0.93 | 0.94 | 346 |

## V. Conclusions

In this paper, we explored the application of random forest to the problem of classifying cars based on their safety features using the Car Evaluation Database. We began by providing an overview of random forest, their structure, and their advantages and limitations. Subsequently, we conducted a comprehensive analysis of the dataset, performed data preprocessing, and trained a random forest classifier. Our findings and key insights from this study can be summarized as follows:

- Random Forests are a powerful and interpretable tool for classification tasks, and their application to car safety classification proved to be effective. - The Car Evaluation Database provided valuable insights into the decision-making process for car safety assessment. - We achieved a respectable level of accuracy in classifying cars into safety categories, demonstrating the feasibility of using random forest for this purpose.

However, it is important to note that our study is not without limitations. Random Forests can be sensitive to noise and may not handle complex relationships between features. Further improvements may be achieved through feature engineering, model hyperparameter tuning, or exploring other machine learning algorithms.

### A. Future Work

There are several avenues for future research in this domain:

- **Feature Engineering:** Exploring additional features or engineering existing ones to improve model performance and predictive accuracy.

- **Comparative Analysis:** Conducting a comparative analysis of Random forest models with other machine learning algorithms to determine the most suitable approach for car safety classification.

- **Real-world Application:** Extending this research to practical applications in the automotive industry, such as automated safety assessment for vehicle manufacturing and insurance.

- **Data Expansion:** Collecting and incorporating more diverse and extensive datasets to enhance the model's generalization and prediction capabilities.

This paper serves as a foundation for future work, and we hope it inspires further exploration of random forest techniques in solving real-world classification problems, particularly in the domain of car safety assessment.

## References

[1] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.

[2] Scikit-learn: Machine Learning in Python. https://scikit-learn.org/stable/

[3] Bohanec,Marko. (1997). Car Evaluation. UCI Machine Learning Repository. https://doi.org/10.24432/C5JP48. https://archive.ics.uci.edu/dataset/19/car+evaluation

[4] Waskom, M. (2021). seaborn: Statistical Data Visualization. Retrieved from https://seaborn.pydata.org/

[5] Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. Computing in Science & Engineering, 9(3), 90-95.

[6] McKinney, W. (2010). Data Structures for Statistical Computing in Python. Proceedings of the 9th Python in Science Conference, 51-56.

[7] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: with Applications in R*, Springer, 2013.

[8] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, 2009.

[9] W. McKinney, *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython*, 2018.

[10] A. Smith, *Advanced Data Analysis Techniques*, Publisher, 2020.

[11] B. Jones, *Statistical Methods for Data Science*, Publisher, 2019.