

# Random Forest Investigation

Freddy Alonzo Mondragon, Gabriela Elizabeth Avila Chan, Mariana Chi Centeno, Karen Valeria Villanueva Novelo, Manuel Nicolas Nahib Franco Valencia, Diego Armando May Pech *Computational Robotics Engineering, Universidad Politecnica de Yucatan*

**Abstract**—Random Forest, an ensemble learning method, stands as a robust and highly accurate approach in the realm of machine learning. It harnesses the collective wisdom of multiple decision trees to deliver predictions that are both dependable and resilient. This versatile technique boasts several advantages, including superior accuracy, effective reduction of overfitting, and the ability to effortlessly accommodate both categorical and numerical data. It unveils feature importance, shedding light on the critical factors driving its predictions, and demonstrates an exceptional tolerance for outliers. Scaling gracefully to large datasets and harnessing the power of parallelization, Random Forest proves efficient and adaptable for various tasks.

Yet, the path to excellence is not without its challenges. As the number of trees increases, the model's complexity and interpretability may diminish. Training may require more time compared to individual decision trees, and the model's inner workings often remain concealed behind the "black box" label. Resource consumption can rise with an abundance of trees, and a certain bias towards dominant classes can surface in imbalanced datasets.

**Index Terms**—Random Forest, Machine Learning, Classification, Regression, Interpretability, Data Preprocessing, Non-linear Relationships, Feature Selection, Scalability, Outliers

## I. INTRODUCTION

THE random forest algorithm, an ensemble learning technique, is a cornerstone in the machine learning landscape, renowned for its versatility and reliability. It operates by constructing a multitude of decision trees at training time and outputting the mode of the classes (classification) or mean prediction (regression) of the individual trees. This method has been successfully applied in various domains ranging from finance and healthcare to environmental science, where predictive accuracy and robustness against overfitting are paramount.

One of the primary merits of random forests lies in their intrinsic nature of combining the simplicity of decision trees with the power of ensemble learning to enhance predictive performance. Each tree in the forest is built on a random subset of data and a random selection of features, leading to diverse trees that collectively contribute to a more stable and accurate model. The aggregation of predictions from multiple trees—be it through voting for classification or averaging for regression—tends to balance out biases, reduce variance, and improve the model's generalization capabilities on unseen data.

Moreover, random forests are remarkably user-friendly, requiring very little hyper-parameter tuning to achieve satisfactory results, which contrasts with more sensitive

algorithms that necessitate careful optimization. This ease of use, combined with their robust performance metrics, has cemented random forests as a go-to algorithm for practitioners looking for an effective off-the-shelf predictive model. As such, they continue to be a subject of extensive research, with ongoing improvements and adaptations expanding their applicability and efficiency in tackling complex machine learning challenges.

## II. HOW RANDOM FOREST WORKS:

Random forest is a commonly-used machine learning algorithm trademarked by Leo Breiman and Adele Cutler, which combines the output of multiple decision trees to reach a single result. Its ease of use and flexibility have fueled its adoption, as it handles both classification and regression problems.

Random forest algorithms have three main hyperparameters, which need to be set before training. These include node size, the number of trees, and the number of features sampled. From there, the random forest classifier can be used to solve for regression or classification problems. Here's an explanation of how Random Forest works:

- Step 1: Select random samples from a given data or training set.
- Step 2: This algorithm will construct a decision tree for every training data.
- Step 3: Voting will take place by averaging the decision tree.
- Step 4: Finally, select the most voted prediction result as the final prediction result.

This combination of multiple models is called Ensemble. Ensemble uses two methods:

- Bagging: Creating a different training subset from sample training data with replacement is called Bagging. The final output is based on majority voting.
- Boosting: Combining weak learners into strong learners by creating sequential models such that the final model has the highest accuracy is called Boosting. Example: ADA BOOST, XG BOOST.

Bagging: From the principle mentioned above, we can understand Random forest uses the Bagging code. Now, let us understand this concept in detail. Bagging is also known

as Bootstrap Aggregation used by random forest. The process begins with any original random data. After arranging, it is organised into samples known as Bootstrap Sample. This process is known as Bootstrapping. Further, the models are trained individually, yielding different results known as Aggregation. In the last step, all the results are combined, and the generated output is based on majority voting. This step is known as Bagging and is done using an Ensemble Classifier.

#### A. Mathematics Behind Random Forest:

The main mathematics behind Random Forest involves understanding how the ensemble combines predictions and why it works effectively.

- **Bootstrap Sampling:** The probability of an instance being selected for a bootstrap sample is

$$1 - \left(1 - \frac{1}{n}\right)^n \quad (1)$$

where  $n$  is the number of instances in the dataset. This probability is derived from the complementary probability of not selecting an instance in  $n$  trials.

- **Feature Selection:** For feature selection, Random Forest randomly selects a subset of  $m$  features out of the total  $d$  to consider at each split. The size of this subset is often controlled by a parameter, such as  $\sqrt{d}$  where  $d$  is the total number of features. This randomness in feature selection contributes to the diversity of the trees and helps in reducing over-fitting.
- **Voting (Classification) or Averaging (Regression):** In classification tasks, the Random Forest ensemble chooses the class with the most votes. In regression, it averages the predictions made by individual trees.

$$\hat{y} = \arg \max_i \sum_{k=1}^K v_{ki} \quad (2)$$

where  $v_{ki}$  is the vote of the  $k$ -th tree for class  $i$ . For regression, the final prediction is the average of the predictions:

$$\hat{y} = \frac{1}{K} \sum_{k=1}^K y_k \quad (3)$$

where  $y_k$  is the prediction by the  $k$ -th tree.

- **Reduction in Over-fitting:** By aggregating the predictions of multiple decision trees, each trained on different subsets of the data and features, Random Forest mitigates the risk of over-fitting. This benefit is attributed

to the Law of Large Numbers, which states that the average of the results obtained from a large number of trials should be close to the expected value, and will tend to become closer as more trials are performed.

The exact mathematical details for optimizing feature selection and building decision trees within a Random Forest are complex and rely on impurity measures (similar to decision trees) and information gain or Gini index for selecting the best features and splits. These details can get quite involved and are usually handled by machine learning libraries like scikit-learn.

In practice, the mathematics behind Random Forest emphasizes the principles of ensemble learning, bootstrap sampling, and the power of averaging or voting to make robust predictions. The combination of these principles is the key to its success in various machine learning tasks.

### III. IMPLEMENTATION OF THE CODE

#### A. Example of Random Forest Code

```

1 # Import necessary libraries
2 from sklearn.datasets import load_iris
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score
6
7 # Load the Iris dataset
8 iris = load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the dataset into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(
14     X, y, test_size=0.3, random_state=42)
15
16 # Create a Random Forest Classifier
17 clf = RandomForestClassifier(n_estimators=100,
18     random_state=42)
19
20 # Train the classifier on the training data
21 clf.fit(X_train, y_train)
22
23 # Make predictions on the test data
24 y_pred = clf.predict(X_test)
25
26 # Evaluate the accuracy of the model
27 accuracy = accuracy_score(y_test, y_pred)
28 print(f"Accuracy: {accuracy*100:.2f}%")

```

Listing 1. Example of Random Forest

#### B. Code Explanation:

```

1 from sklearn.datasets import load_iris
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score

```

The code begins by importing the necessary libraries from scikit-learn, a popular machine learning library in Python.

```
1 iris = load_iris()
2 X = iris.data
3 y = iris.target
```

The code uses the `load_iris()` function to load the Iris dataset, which is a well-known dataset used for classification tasks. It contains measurements of sepal and petal lengths and widths for three different species of iris flowers. The feature data (sepal and petal measurements) is stored in the variable `X`, and the target labels (species of iris) are stored in the variable `y`.

```
1 X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42)
```

The dataset is divided into training and testing sets using the `train_test_split` function. In this case, 70% of the data is used for training (`X_train` and `y_train`), and 30% is used for testing (`X_test` and `y_test`). The `random_state` parameter is set to 42 to ensure reproducibility in the random splitting.

```
1 clf = RandomForestClassifier(n_estimators=100,
    random_state=42)
```

A Random Forest Classifier is created using the `RandomForestClassifier` class from `scikit-learn`. `n_estimators=100` specifies that the Random Forest will consist of 100 decision trees. You can adjust this number based on your needs. The `random_state` parameter is set for reproducibility.

```
1 clf.fit(X_train, y_train)
```

The `fit` method is used to train the Random Forest Classifier on the training data (`X_train` and `y_train`).

```
1 y_pred = clf.predict(X_test)
```

After training, the model is used to make predictions on the test data (`X_test`), and the predicted labels are stored in the variable `y_pred`.

```
1 accuracy = accuracy_score(y_test, y_pred)
2 print(f"Accuracy: {accuracy*100:.2f}%")
```

The `accuracy_score` function is used to calculate the accuracy of the model by comparing the predicted labels (`y_pred`) with the true labels from the testing set (`y_test`). The code prints the accuracy in percentage form.

The significance of this code in the context of the topic at hand (Random Forest) lies in its demonstration of how to use Random Forest for a classification task. The result obtained, the model's accuracy, is pivotal in assessing how well the model can make predictions on new data. In this case, a high accuracy would indicate that the Random Forest model is effective in classifying Iris flower species in the test data-set. Accuracy is a critical metric for measuring the performance of classification models and serves as an indicator of their ability to generalize and make precise predictions in real-world scenarios.

When you execute the code, you will obtain an accuracy value (as a percentage) printed in the console, which will provide you with insights into how well the model performed on the Iris data-set's test data. For instance, if you achieve an accuracy of 0.95, it means the model correctly classified 0.95 of the samples in the test data-set, indicating a robust performance.

#### IV. ADVANTAGES AND DISADVANTAGE

##### A. Advantages of Random Forest:

- **High Accuracy:** Random Forest is known for its high predictive accuracy. It often outperforms single decision trees and is a go-to choice for many classification and regression tasks.
- **Reduced Overfitting:** The ensemble nature of Random Forest, where multiple decision trees are combined, helps mitigate overfitting. By averaging or voting on multiple trees, the model generalizes well to new, unseen data.
- **Feature Importance:** Random Forest provides a feature importance score, allowing you to identify which features have the most significant impact on the model's predictions. This can aid in feature selection and data understanding.
- **Handles Both Categorical and Numerical Data:** Random Forest can work with a mix of categorical and numerical data without requiring extensive preprocessing. It automatically handles missing values.
- **Robust to Outliers:** Random Forest is relatively robust to outliers in the data. Outliers may affect individual decision trees but have less impact when combined in an ensemble.
- **Scalability:** It can handle large datasets and high-dimensional feature spaces effectively.
- **Parallelization:** Training Random Forest can be parallelized, allowing for faster training on multi-core processors.
- **Implicit Feature Selection:** The algorithm can implicitly perform feature selection by giving higher importance to the most informative features when making splits in decision trees.

### B. Disadvantages of Random Forest:

- **Complexity:** Random Forest models can become complex, especially when using a large number of trees. Interpreting and visualizing these models can be challenging.
- **Slower Training:** Random Forest can be slower to train compared to single decision trees due to the ensemble nature. Training more trees in the forest takes more time.
- **Lack of Transparency:** While Random Forest is highly accurate, it is often considered a "black box" model because it's challenging to understand the individual decision-making process of each tree in the forest.
- **Resource Intensive:** In terms of memory and computation, Random Forest can be resource-intensive when dealing with a large number of trees.
- **Bias Toward Dominant Classes:** In imbalanced datasets, Random Forest can be biased toward the majority class. This means it might not perform as well on minority classes.
- **Limited Performance Improvement:** After a certain point, adding more trees to the forest may not significantly improve performance but can increase computational costs.

## V. FIELDS OF APPLICATION FOR RANDOM FOREST ALGORITHM:

The Random Forest algorithm finds applications in a myriad of diverse fields, showcasing its versatility in addressing classification and regression tasks. It is recognized for its proficiency in managing an array of data types and mitigating overfitting.

- 1) **Health Sciences:**
  - Medical diagnosis, such as disease detection.
  - Medical image analysis.
  - Prediction of medical treatment outcomes.
- 2) **Finance:**
  - Credit risk assessment.
  - Financial fraud detection.
  - Market movement and stock price prediction.
- 3) **Biology and Genomics:**
  - Species classification in biology.
  - Gene and protein prediction.
  - DNA sequencing data analysis.
- 4) **Agriculture:**
  - Crop yield prediction.
  - Plant disease detection.
  - Soil classification.

### 5) Marketing and Advertising:

- Customer segmentation.
- Product recommendation.
- Sentiment analysis on social media.

### 6) Manufacturing and Quality:

- Quality control in production.
- Predictive maintenance of machinery.
- Optimization of manufacturing processes.

### 7) Human Resources:

- Candidate selection and human resources.
- Employee turnover prediction.
- Employee satisfaction analysis.

### 8) Environment:

- Air and water quality monitoring and prediction.
- Natural disaster prediction.
- Biodiversity analysis.

### 9) Education:

- Student assessment and classification.
- Personalized learning.
- Graduation rate prediction.

### 10) Transportation and Logistics:

- Route and logistics planning.
- Predictive vehicle maintenance.
- Transportation demand prediction.

### 11) Energy:

- Energy demand prediction.
- Predictive maintenance of energy infrastructure.
- Time series analysis in the energy sector.

The Random Forest algorithms broad spectrum of applicability makes it a valuable asset in the realm of data-driven decision-making, addressing the unique demands of each of these domains.

## VI. CONCLUSION

Random Forest is a powerful ensemble learning method that combines the strengths of multiple decision trees to deliver robust and highly accurate predictions. This versatile machine learning technique offers a wealth of advantages, but it's not without its limitations.

The advantages of Random Forest include its remarkable accuracy, reduced overfitting, and the ability to handle both categorical and numerical data. It excels at feature importance analysis, providing valuable insights into the factors driving its predictions. Furthermore, Random Forest is known for its resilience to outliers, scalability to large datasets, and the ability to implicitly perform feature selection. Its parallelization capabilities make it efficient for training on multi-core processors.

However, there are some downsides to consider. Random Forest models can become complex and challenging to interpret, especially when dealing with a large number of

trees. Training time may be slower compared to single decision trees due to the ensemble approach. While the model offers high accuracy, it is often considered a "black box," lacking transparency in individual tree decision-making. Additionally, Random Forest can be resource-intensive, especially when working with numerous trees, and it may exhibit bias toward dominant classes in imbalanced datasets.

## REFERENCES

- [1] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [2] T. K. Ho, "Random decision forests," *Proceedings of the 3rd International Conference on Document Analysis and Recognition*, pp. 278-282, 1995.
- [3] T. Hastie, R. Tibshirani, and J. H. Friedman, "The elements of statistical learning," 2nd ed., New York: Springer, 2009.
- [4] P. Geurts, D. Ernst, and L. Wehenkel, "Extremely randomized trees," *Machine Learning*, vol. 63, no. 1, pp. 3-42, 2006.
- [5] A. Liaw and M. Wiener, "Classification and regression by random forest," *R News*, vol. 2, no. 3, pp. 18-22, 2002.
- [6] IBM, "What is random forest? — IBM," [Online]. Available: <https://www.ibm.com/topics/random-forest#:~:text=Random%20forest%20is%20a%20commonly,both%20classification%20and%20regression%20problems.> [Accessed November 7, 2023].
- [7] Simplilearn, "Random forest algorithm," [Online]. Available: <https://www.simplilearn.com/tutorials/machine-learning-tutorial/random-forest-algorithm> [Accessed November 7, 2023].
- [8] "Built In, 'Random forest: A complete guide for machine learning,' [Online]. Available: <https://builtin.com/data-science/random-forest-algorithm#how> [Accessed November 7, 2023]."
- [9] [Interactive Chaos. (s.f.). 'Random Forest.' [Online] <https://interactivechaos.com/es/wiki/random20Forest>