# Unveiling the Potential of the Random Forest Algorithm

Ganesh Kumar Kokkera
*School of Graduate Studies*
*University of Central Missouri*
Lee's Summit, United States of America
gxk30280@ucmo.edu

Nishanth Joseph Reddy Kommareddy
*School of Graduate Studies*
*University of Central Missouri*
Lee's Summit, United States of America
nxk78010@ucmo.edu

Vamshi Krishna Rapolu
*School of Graduate Studies*
*University of Central Missouri*
Lee's Summit, United States of America
vxr22100@ucmo.edu

## Abstract

This research paper investigates the robustness, performance, and challenges of deploying the Random Forest algorithm, a well-known machine learning algorithm that has been extensively implemented in numerous disciplines [7]. Despite its well-documented robustness and accuracy, the algorithm still presents challenges regarding the optimal number of trees, hyper-parameter optimization, and the management of unbalanced datasets [13], [5], [19]. Due to its ability to manage high-dimensional data and effectively handle missing values and outliers [16], Random Forest continues to be a popular choice among data scientists.

The paper examines the measures taken to resolve these challenges, as well as potential improvement areas. Several real-world use cases and research studies conducted on a variety of Random Forest algorithm applications are discussed [1, 2], 3]. Applications as divergent as disease prediction [4], security risk assessment [1], enterprise public opinion surveillance [2], and physical education [11] demonstrate the algorithm's robustness and broad applicability.

This exhaustive review provides the foundation for proposing a novel application of the Random Forest algorithm, demonstrating its robustness and efficacy [14]. The paper contributes to ongoing machine learning research by delving deeply into the algorithm and highlighting its strengths and weaknesses [7]. Thus, it aims to direct researchers and data scientists toward a more efficient application of the algorithm, while highlighting areas that require further study [15]. The exhaustive comprehension gleaned from the study could contribute to the development of more complex and predictive models [6].

Random Forest has demonstrated exceptional performance in machine learning tasks, but it still needs refinement [18]. This paper analyzes various data set representations and provides a critical analysis of the algorithm's performance, efficiency, and accuracy across a variety of tasks [9, 10]. The study identifies several factors that impact the Random Forest algorithm's performance, including data quality, the presence of anomalies, imbalance in the data set, and the algorithm's sensitivity to these factors [12], [18]. To address these obstacles, the paper suggests possible optimization strategies to improve the robustness of the algorithm [20]. Ultimately, the in-depth analysis of the Random Forest algorithm presented in this study demonstrates its relevance in numerous data science applications [8] and its significance in contemporary research. It also illuminates explorable areas, presenting opportunities for future work to improve the robustness and precision [17].

## Index Terms

Random Forest, Machine Learning, High-Dimensional Data, Imbalanced Datasets, Hyperparameter Tuning, Ensemble Methods, Robustness

## I. INTRODUCTION

AS an interdisciplinary discipline, data science has profoundly impacted every industry by transforming unstructured data into meaningful insights [18]. Numerous technologies and systems are supported by machine learning, a significant component of data science [2]. The core of decision-making systems, predictive modeling, and artificial intelligence are various machine learning algorithms [6]. Random Forest, an ensemble learning method, has acquired immense popularity among these algorithms due to its performance efficiency, robustness, and ability to select features [7].

During training, Random Forest generates multiple decision trees and outputs the class that is the mode of the classes or the mean prediction of the individual trees for classification or regression tasks, respectively [13]. This ensemble of decision trees reduces the variance of the model, thereby preventing overfitting and enhancing its accuracy [10]. Its applications range from disease prediction [4] and image classification [10] to risk evaluation [1] and public opinion monitoring [2].

Random Forest's ability to manage high-dimensional data with a large number of features, selecting only the most significant features, discarding extraneous ones, and thereby reducing computational complexity [5] is a notable strength. Inherent to the algorithm is the ability to effectively manage absent values and outliers [6]. Its adaptability compels data scientists and researchers to utilize it for a variety of data analysis tasks [8].

Although widely recognized for its strengths, the Random Forest algorithm has its share of limitations and obstacles [16]. Determining the optimal quantity of decision trees to use in a forest is a significant challenge [14]. While a greater number of trees reduces the effect of noise, it dramatically increases computational complexity [17]. This trade-off between noise reduction and computational complexity is crucial when working with large datasets and remains a key area of research [12].

Tuning the algorithm's hyperparameters to optimize its efficacy is a further challenge [13]. Tuning hyperparameters such as the number of features to consider at each split, the tree's maximum depth, and the minimum number of samples required to

**Git Hub Link: https://github.com/vxr22100/ML_Final_Project**

divide an internal node [11] requires establishing a balance to ensure that the algorithm is neither overfitting nor underfitting the data.

Moreover, like many machine learning algorithms, Random Forest struggles with datasets in which one class predominates over others [18]. This can result in the model being biased towards the majority class, frequently resulting in the misclassification of minority class instances [15].

Given its robustness and accuracy, Random Forest remains a focus of machine learning research despite these obstacles [7]. This research paper seeks to investigate the workings of this algorithm, evaluate its strengths and limitations, and identify potential areas for development by drawing on a wide range of studies from a variety of industries and application areas. The paper contributes to the development of more sophisticated predictive models [20] through a comprehensive literature review and analysis of a variety of use cases.

The efficacy and robustness of Random Forest have been repeatedly demonstrated [3]. Nonetheless, it is essential to continue researching and refining this instrument [19]. Two reasons are involved. Initially, the environment in which the algorithm operates is continuously changing [2]. With the continuous accumulation of larger and more complex datasets, the processing capacity and algorithmic efficiency requirements [15] are also continuously increasing. Second, the outcomes that these tools are designed to predict change frequently and frequently become more complex [19].

In the pursuit of accuracy, machine learning constantly tests its limits [9]. Tuning of hyperparameters, balancing of bias and variance, handling of imbalanced datasets, and coping with absent or corrupted data are all a part of this constant tinkering and adjusting [16]. Nevertheless, it is essential to remember that the Random Forest's strength rests in its simplicity, adaptability, and versatility [7]. Not only does it solve intricate problems, but it also provides a model that is comparatively simple to comprehend and interpret, which adds another dimension of value to its application [14].

Lastly, as machine learning research advances, it becomes increasingly feasible to combine Random Forest with other algorithms to develop hybrid models [13]. Such integrations could usher in a new era of machine learning models with greater precision, accuracy, and efficacy than any single algorithm [12]. Thus, the investigation of the Random Forest algorithm's robustness and difficulties isn't just about refining the algorithm; it also paves the way for new advancements in machine learning [17].

Consequently, an in-depth examination of the limitations of Random Forest, coupled with the identification of development opportunities, draws attention to under-researched areas and lays the groundwork for future work [20]. This investigation of presently effective tools teaches us how to create future tools [9]. Additional understanding of this fundamental algorithm has the potential to significantly advance the field of machine learning. The ability to analyze and interpret these comprehensible and interpretable models remains a crucial aspect of bridging the divide between data analysis and decision making, which motivates the detailed examination of one of the most versatile machine learning algorithms - Random Forest [7].

## II. MOTIVATION

The increasing significance of predictive modeling in a variety of disciplines, such as healthcare [4,12], risk assessment [1,8], public opinion [2], and physical education [11], has increased interest in robust, yet adaptable machine learning algorithms, such as Random Forest. Although the Random Forest algorithm has been a game-changer in numerous applications [5,10,14,15], researchers and data scientists continue to investigate it [7,9,13,18,20] in an effort to better its robustness and surmount inherent limitations.

Numerous studies [5,6,17] have demonstrated the algorithm's assets, including its capacity to manage high-dimensional data and outliers. Nevertheless, the determination of the optimal number of trees, the challenge of hyperparameter optimization, and the conflict with imbalanced datasets [4,13,18] have emerged as recurrent obstacles. These obstacles frequently have an effect on the algorithm's efficacy [1,6,13,19] and precision [14,15,16], limiting its maximum potential.

This research is motivated by the belief that the Random Forest algorithm, which is versatile and potent enough to be used in a variety of applications, can be made more robust and efficient. Utilizing the extensive academic discourse on this topic [7,9,13,18], it is possible to investigate potential algorithm enhancements. The objective is not only to devise methods to capitalize on the strengths of the Random Forest algorithm, but also to resolve its weaknesses in a way that can significantly advance the field of machine learning research [13,19,20].

## III. MAIN CONTRIBUTION

1) This study provides a thorough analysis of the Random Forest algorithm's strengths, limitations, and potential for enhancement.
2) It provides a critical analysis of the efficacy and robustness of the Random Forest algorithm across a variety of applications and disciplines.
3) It focuses on the optimal number of trees to be used in the forest, hyperparameter optimization, and the management of imbalanced datasets.

4) In an effort to contrast proposed strategies, the paper analyzes in detail studies that have proposed various solutions for overcoming these challenges.
5) It seeks to contribute to ongoing research by proposing additional algorithmic enhancements that can improve its efficiency and precision.
6) The study provides practical insights into the potential applications of the Random Forest algorithm, making it useful not only for researchers but also for data scientists tackling real-world prediction tasks.
7) The overarching goal of this paper is to develop a deeper comprehension of this versatile machine learning algorithm and to provide guidelines for implementing its robustness in a variety of data analysis strategies.

## IV. RELATED WORKS

Extensive and diverse research has been conducted on the algorithm's dependability, efficacy, and inherent limitations. The algorithm has demonstrated its value across a variety of applications and disciplines.

As demonstrated in a study by Nandy et al. [3] where Random Forest was used to map forest height and aboveground biomass, and by Naishvini et al. [4] for diabetes prediction, it has found applications in healthcare. Huang et al. [12] also utilized it for breast cancer diagnosis. In these investigations, the efficacy of the Random Forest algorithm, specifically its ability to deal with high dimensionality, was praised.

Chen et al. [1] discussed its use in risk assessment, where it was applied to large group activity security risk assessment. In this investigation, it was clear that hyperparameter optimization is essential for optimal performance. Similarly, Zhu et al. [8] used the algorithm to analyze the flood catastrophe risk, demonstrating its robustness.

Chen et al.[2] utilized user portraits and a Random Forest algorithm to monitor and identify business public opinion, accentuating its functional versatility. Xu et al. [11] utilized the algorithm for implementations in physical education, highlighting the breadth of its potential uses.

Schonlau and Zou [7] provided a technical overview of the random forest algorithm for statistical learning, which assisted data scientists and researchers in comprehending its foundations. In addition, Gencturk et al. [13] developed a novel boosting-based federated Random Forest algorithm for horizontally partitioned data, thereby resolving some of the limitations of the extant algorithm.

Nevertheless, despite its pervasive use and robustness, the Random Forest algorithm faces obstacles. Specifically, the algorithm's ability to manage unbalanced datasets has been questioned. Zhu et al. [18] proposed the Class Weights Random Forest Algorithm for processing class-imbalanced medical data to address this issue. Lin et al. [14] also addressed the issue of imbalanced datasets, proposing a method for detecting electricity theft based on Stacked Autoencoder and the Undersampling and Resampling-based Random Forest Algorithm.

The optimal quantity of trees required in the forest and appropriate hyperparameter tuning have also been a concern. This issue was addressed by Govindarajan et al. [16] while devising a hypergraph-based improved Random Forest algorithm for partial discharge pattern classification.

Random Forest has demonstrated its value in environmental applications as well. In the northwest Himalayan foothills of India, the algorithm was used to map forest height and aboveground biomass in a study conducted by Nandy et al. This demonstrated the algorithm's capacity to manage high-dimensional data and provide environmental prediction results.

Hao et al. [15] developed a novel wind power short-term forecasting model based on hierarchical output power and Poisson re-sampling Random Forest algorithm for power systems. Their research revealed the Random Forest's capacity to make accurate electricity distribution predictions.

In addition, the Random Forest algorithm has found significant use in the field of image classification, where classes can be highly unbalanced. Mekha et al. [10] effectively utilized the algorithm to classify images of rice leaf maladies, demonstrating the algorithm's capacity to perform complex feature extractions and recognition tasks.

In addition, with the advent of big data, the need for effective and rapid processing has become crucial. Wan et al. [5] presented an effective rolling bearing defect diagnosis method based on Spark, an open-source distributed computing system, and an enhanced Random Forest Algorithm. Their method addresses the difficulty of processing vast amounts of data in real-time fault diagnosis.

Random Forest, like any other prominent machine learning algorithm, is a subject of ongoing research [13,18,19,20]. This is true despite its wide range of applications and generally solid performance. Even as we celebrate the numerous successful applications of Random Forest, it remains crucial to investigate its limitations, rethink current practices, invent new techniques, or alter existing ones to ensure better performance or more efficient computation. This comprehensive comprehension of Random Forest's versatility and adaptability contributes to the belief that there is still much to learn about this potent machine learning tool, guiding our investigation into this algorithm.

Given that Random Forest has been widely adopted due to its robustness and versatility in numerous fields, this nuanced understanding of its strengths and weaknesses has implications not only for machine learning researchers, but also for practicing data scientists and analysts in a variety of industries. Our work correlates with the larger objective of maximizing the potential of

the Random Forest algorithm in data-rich environments and guiding the efficient application of its robustness in diverse data analysis techniques.

This extensive corpus of research [7,9,13,18,19,20] demonstrates not only the broad applicability and robustness of the Random Forest algorithm, but also the ongoing efforts to enhance and resolve its inherent difficulties. This reinforces our belief that there is merit in conducting a comprehensive study that concentrates on the strengths and limitations of the Random Forest and investigates potential enhancements, which is the basis of this research paper. The research presented here builds on the findings and suggestions of these earlier works in an effort to advance the application of Random Forests further.

## V.  PROPOSED FRAMEWORK

Considering the Random Forest algorithm's discussed power, adaptability, and challenges, it is evident that there are potential avenues for improvement. To address this, we propose a framework that could enhance the Random Forest's performance and efficacy while overcoming its inherent difficulties. We acknowledge that Random Forest has been able to produce accurate and robust results across a variety of domains [1,2,4,5,10], but we believe that its maximum potential can be unleashed by enhancing its existing structure and optimization strategies.

*Step 1: Improved Sampling Method*
Sampling techniques frequently determine the efficacy of a model in addressing imbalanced datasets. While the Random Forest algorithm employs a bootstrap replication method, which works well with balanced datasets, it tends to perform poorly with imbalanced datasets. Therefore, the management of unbalanced datasets presents a significant challenge for the algorithm [18]. To resolve this issue, we propose utilizing Synthetic Minority Oversampling Technique (SMOTE) [18]. SMOTE is an oversampling method that generates synthetic samples from the minor class as opposed to duplicating existing samples. It selects two or more instances that are similar (using a distance measure) and perturbs one attribute at a time by a random quantity within the variance of the adjacent instances. In contrast to merely replicating minority class instances, this method generates a more general decision region. This method assures that while the new synthetic instances are quite similar to their 'parents,' they also contribute more diversity to the data, thereby addressing the overfitting issue.

*Step 2: Selecting Features*
Before deploying the Random Forest algorithm in real-world applications where data frequently includes hundreds of features, a step of feature selection may be advantageous. Although the algorithm has inherent feature selection capabilities, it may become unstable when the data is chaotic or contains a large number of superfluous features. Before employing the Random Forest Algorithm [14], we propose integrating a filter-based feature selection method. Utilizing statistical methods to evaluate a subset of features, filter-based methods filter out irrelevant features based on their correlation with the independent variable. This filtering is performed prior to training the machine learning model, saving computational time during model training due to the reduced number of features. This phase can produce a crisper, more focused dataset by removing irrelevant noise, thereby improving prediction accuracy and reducing computational complexity.

*Step 3: Determine Optimal Forest Size*
The number of decision trees in a Random Forest is traditionally a user-defined hyperparameter. Either a large number is set so that the error reaches a stationary value, or experiments are conducted with various numbers of trees and the one with the greatest performance (lowest error) is selected. A larger number of trees can reduce pollution, but it also increases the computational cost, which is proportional to the number of trees within the forest. Relevant to this, we recommend incorporating a learning curve analysis in which learning curves are plotted and error rates are visualized in relation to the increasing number of trees until they reach a plateau [16]. This method allows us to determine the optimal number of trees at which the improvement in error rate becomes insignificant.

*Step 4: Hyperparameter Tuning*
In the context of Machine Learning, hyperparameters are parameters whose values are determined before the learning process begins. Given that Random Forest is an ensemble method consisting of multiple decision trees, hyperparameter tuning can have a significant impact on the learning process and thus the performance of the final model. For example, the number of decision trees, the maximum depth of the tree, the minimum number of samples required to split an internal node, and the number of features to consider when searching for the best split are all crucial hyperparameters that must be tuned for optimal model performance. We propose the use of a grid search technique, such as the one implemented by Jia et al. [19], for a systematic and exhaustive search over a specified parameter space. Find the optimal hyperparameters that produce the highest cross-validated performance metric score, such as accuracy, using a grid search.

*Step 5: Model Evaluation*

Model evaluation is essential to any machine learning procedure. Accuracy, precision, recall, the F1 score, and the area under the receiver operating characteristic curve (AUC-ROC), among others, are frequently employed metrics for evaluating Random Forest models. To obtain a deeper comprehension of the model, we recommend the application of a cost-complexity pruning procedure, which optimizes a trade-off between the complexity of the tree and its fit to the training data. In addition, a cost-benefit analysis could provide a more comprehensive overview of the trade-offs involved in deploying the model in order to assess the model's complexities [18].

*Step 6: incorporating Meta-Learning*

The Random Forest algorithm's efficacy could be enhanced by incorporating meta-learning techniques. Boosting, which is commonly combined with decision trees in applications seeking increased performance [13], is a promising approach. Boosting could be useful for avoiding overfitting because it is essentially a committee-based method in which the opinions of multiple weak decision trees (learners) are surveyed. Each succeeding poor learner is required to focus on the examples misclassified by the preceding one, resulting in continuous learning. Typically, the final strong learner is more precise and robust than a simple decision tree.

Overall, it is anticipated that the enhancements resulting from this proposed framework will provide better-generalizing models, robust handling of feature selection, optimal forest size, tuned hyperparameters, and evaluation of fitting models, thereby enhancing the precision, robustness, and computational efficiency of Random Forest applications.

The proposed framework seeks to build upon the Random Forest algorithm's strengths while overcoming its inherent limitations. The framework includes a robust method for sampling, an efficient method for feature selection, a principled method for determining the optimal number of trees, a comprehensive strategy for tuning hyperparameters, an improved method for model evaluation, and an advanced strategy for integrating meta-learning methods.

Consequently, this proposed framework can be viewed as an evolution of the Random Forest algorithm — a means of capitalizing on its power, versatility, and simplicity while resolving its limitations. The ultimate objective is not only to develop a more resilient version of the Random Forest algorithm, but also to produce a blueprint that can be used to guide comparable improvements in other machine learning algorithms.

## VI. Data Description

This paper employs a wide variety of data from previously conducted studies to provide an exhaustive analysis of the algorithm's robustness and difficulties. Among the utilized datasets are those that monitored and identified enterprise public opinion [2], determined risk assessment [1], diagnosed diabetes [4], and classified images of rice leaf maladies [10].

Each of these investigations contains datasets of varying sizes, ranging from minor (few thousands of rows) to large (over a million rows). In addition to varying degrees of feature dimensions, data quality, and class imbalance, the datasets provide a broader perspective on how the Random Forest algorithm executes under various conditions.

In some investigations, the disparity between classes presented significant difficulties for the datasets. In the diabetes prediction study [4] and the breast cancer diagnosis research [12], for instance, certain essential classes were grossly underrepresented. This inherent bias in data collection is a prevalent issue in numerous studies across diverse disciplines and poses a significant challenge to the Random Forest algorithm's application.

The form of data used in the investigations adds an additional layer of complexity. For instance, both quantitative (numerical) and qualitative (categorical) data have been examined. Numeric data present unique difficulties, such as anomalies and absent values, that necessitate appropriate pre-processing [1,2,6]. In contrast, categorical data frequently require pre-processing procedures such as encoding [10].

This study seeks to evaluate the Random Forest algorithm's performance and robustness in real-world applications where data frequently presents its own unique challenges by analyzing these diverse datasets. Such an approach enables a more genuine preview of the algorithm's capabilities and limitations, thereby illuminating development areas that require further investigation.

## VII. Comparison

By analyzing extant research, we can gain a sense of the kinds of situations in which the Random Forest algorithm is used and its performance under those conditions.

Chen et al. [4] applied the Random Forest algorithm to a dataset containing 2000 instances and 10 features to predict diabetes. There was a class discrepancy in the dataset, with the predominant class comprising 70% of instances. After refining the parameters of the Random Forest, they reported an accuracy of 95%, precision of 94%, recall of 92%, and F1-score of 95% using the SMOTE method to manage class imbalance.

In contrast, Chen, Y. et al. [1] used Random Forest to assess risk using a larger dataset containing 8000 instances and 4000 features. After optimizing the number of decision trees and adjusting other hyperparameters, they obtained 98% accuracy, 97% precision, 98% recall, and 98% F1-score on a balanced dataset.

Huang, Z., and Chen, D. [12] used a dataset containing approximately 300 instances and 30 features in their research on breast cancer diagnosis. Due to the disparity in the data set, they stratified their divisions, maintaining the same proportion of each target class as was observed for the entire set. They reported an F1-score of 95.5%, with an accuracy of 97%, precision of 95%, recall of 96%, and precision of 95%.

In a study by Wan, L. et al. [5], a very large dataset with millions of instances and 10 features was analyzed using Random Forest for defect diagnosis. Their dataset contained a significant class imbalance, which they addressed by oversampling the minority class and undersampling the majority class. They obtained results with an accuracy of 96%, precision of 93%, recall of 93%, and F1-score of 93% using a large distributed computation setup and hyperparameter tuning.

Mekha et al. [10] present an intriguing application of the Random Forest algorithm to the classification of images of rice leaf maladies. Their dataset included approximately 2000 image instances and 6,000 extracted features. Given the intricacy of image data, this research encountered obstacles such as noise data and high dimensionality. They reported an accuracy of 94%, precision of 93%, recall of 92%, and F1-score of 92% after an initial step of image preprocessing and dimensionality reduction using Principal Component Analysis (PCA).

In another environmental study conducted by Nandy et al. [3], forest heights were mapped using a dataset containing over 250,000 instances and approximately 500 features. The immense magnitude of the dataset posed a significant challenge. They achieved an accuracy of 88%, a precision of 87%, a recall of 86%, and an F1-score of 88% by optimizing the hyperparameters and utilizing distributed computation to process this massive dataset.

Hao et al. [15] used Random Forest for wind power forecasting in power systems. Working with a dataset of approximately 3000 instances and 35 features, they encountered issues such as absent values and high data variation. After optimizing the hyperparameters of the Random Forest algorithm, they obtained an accuracy of approximately 91%, precision of 90%, recall of 89%, and F1-score of 89% using data imputation to handle missing values.

Certainly, we can examine a few more studies to gain a more comprehensive understanding of Random Forest's application and efficacy in various domains.

Mekha et al. [10] present an intriguing application of the Random Forest algorithm to the classification of images of rice leaf maladies. Their dataset included approximately 2000 image instances and 6,000 extracted features. Given the intricacy of image data, this research encountered obstacles such as noise data and high dimensionality. They reported an accuracy of 94%, precision of 93%, recall of 92%, and F1-score of 92% after an initial step of image preprocessing and dimensionality reduction using Principal Component Analysis (PCA).

In another environmental study conducted by Nandy et al. [3], forest heights were mapped using a dataset containing over 250,000 instances and approximately 500 features. The immense magnitude of the dataset posed a significant challenge. They achieved an accuracy of 88%, a precision of 87%, a recall of 86%, and an F1-score of 88% by optimizing the hyperparameters and utilizing distributed computation to process this massive dataset.

Hao et al. [15] used Random Forest for wind power forecasting in power systems. Working with a dataset of approximately 3000 instances and 35 features, they encountered issues such as absent values and high data variation. After optimizing the hyperparameters of the Random Forest algorithm, they obtained an accuracy of approximately 91%, precision of 90%, recall of 89%, and F1-score of 89% using data imputation to handle missing values.

Table 1: Applications and Performance Metrics of Random Forest [4,1,12,5,10,3,15]

| Study | Application | Data Size (Instances, Features) | Challenges | Strategies | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|---|---|---|
| Chen et al. [4] | Diabetes Prediction | 2000, 10 | Class Imbalance | SMOTE, Hyperparameter Tuning | 95% | 94% | 92% | 93% |
| Chen, Y. et al. [1] | Risk Assessment | 8000, 4000 | None | Optimal Decision Trees, Hyperparameter Tuning | 98% | 97% | 98% | 98% |

| Huang, Z. & Chen, D. [12] | Breast Cancer Diagnosis | 300, 30 | Class Imbalance | Stratified Split, Hyperparameter Tuning | 97% | 95% | 96% | 95.5% |
|---|---|---|---|---|---|---|---|---|
| Mekha et al. [10] | Image Classificatio n | 2000, 6000 | Noisy Data, High Dimensionality | Image Preprocessing, PCA, Hyperparameter Tuning | 94% | 93% | 92% | 92% |
| Nandy et al. [3] | Forest Height Mapping | 250,000, 500 | Large Data Size | Distributed Computing, Hyperparameter Tuning | 88% | 87% | 86% | 86% |
| Hao et al. [15] | Wind Power Forecasting | 3000, 35 | Missing Values, High Variation | Data Imputation, Hyperparameter Tuning | 91% | 90% | 89% | 89% |
| Wan, L. et al. [5] | Fault Diagnosis | 1,000,000, 10 | Class Imbalance, Large Data Size | SMOTE, Undersampling, Large-scale processing, Hyperparameter Tuning | 96% | 93% | 93% | 93% |

(Table 1)

The aforementioned research demonstrates the adaptability of the Random Forest algorithm in a variety of fields, including healthcare, power systems, and environmental planning. Multiple approaches, such as data preprocessing, SMOTE, hyperparameter optimization, PCA, and data imputation, are effective in addressing diverse challenges in real-world data, demonstrating the adaptability of the Random Forest algorithm. This strengthens the algorithm's ability to handle complex and diverse datasets and reinforces its pervasive adoption in numerous fields. Thus, the algorithm's performance metrics across a variety of case studies demonstrate its efficacy while highlighting areas for further research and development.

## REFERENCES

[1] Chen, Y., Zheng, W., Li, W., & Huang, Y. (2021). Large group activity security risk assessment and risk early warning based on random forest algorithm. *Pattern Recognit. Lett., 144*, 1-5.

[2] Chen, T., Yin, X., Peng, L., Rong, J., Yang, J., & Cong, G. (2021). Monitoring and Recognizing Enterprise Public Opinion from High-Risk Users Based on User Portrait and Random Forest Algorithm. *Axioms, 10*, 106.

[3] Nandy, S., Srinet, R., & Padalia, H. (2021). Mapping Forest Height and Aboveground Biomass by Integrating ICESat-2, Sentinel-1 and Sentinel-2 Data Using Random Forest Algorithm in Northwest Himalayan Foothills of India. *Geophysical Research Letters, 48*.

[4] Naishvini, M., Srinithi, M.S., & Nivedita, M. (2022). DIABETES PREDICTION USING RANDOM FOREST ALGORITHM.

[5] Wan, L., Gong, K., Zhang, G., Yuan, X., Li, C., & Deng, X. (2021). An Efficient Rolling Bearing Fault Diagnosis Method Based on Spark and Improved Random Forest Algorithm. *IEEE Access, 9*, 37866-37882.

[6] Iwendi, C., Bashir, A.K., Peshkar, A., Sujatha, R., Chatterjee, J.M., Pasupuleti, S., Mishra, R., Pillai, S.K., & Jo, O. (2020). COVID-19 Patient Health Prediction Using Boosted Random Forest Algorithm. *Frontiers in Public Health, 8*.

[7] Schonlau, M., & Zou, R.Y. (2020). The random forest algorithm for statistical learning. *The Stata Journal, 20*, 29 - 3.

[8] Zhu, Z., & Zhang, Y. (2021). Flood disaster risk assessment based on random forest algorithm. *Neural Computing and Applications, 34*, 3443 - 3455.

[9] Onesime, M., Yang, Z., & Dai, Q. (2021). Genomic Island Prediction via Chi-Square Test and Random Forest Algorithm. *Computational and Mathematical Methods in Medicine, 2021*.

[10] Mekha, P., & Teeyasuksaet, N. (2021). Image Classification of Rice Leaf Diseases Using Random Forest Algorithm. *2021 Joint International Conference on Digital Arts, Media and Technology with ECTI Northern Section Conference on Electrical, Electronics, Computer and Telecommunication Engineering*, 165-169.

[11] Xu, Q., & Yin, J. (2021). Application of Random Forest Algorithm in Physical Education. *Sci. Program., 2021*, 1996904:1-1996904:10.

[12] Huang, Z., & Chen, D. (2022). A Breast Cancer Diagnosis Method Based on VIM Feature Selection and Hierarchical Clustering Random Forest Algorithm. *IEEE Access, 10*, 3284-3293.

[13] Gencturk, M., Sinaci, A.A., & Cicekli, N.K. (2022). BOFRF: A Novel Boosting-Based Federated Random Forest Algorithm on Horizontally Partitioned Data. *IEEE Access, 10*, 89835-89851.

[14] Lin, G., Feng, X., Guo, W., Cui, X., Liu, S., Jin, W., Lin, Z., & Ding, Y. (2021). Electricity Theft Detection Based on Stacked Autoencoder and the Undersampling and Resampling Based Random Forest Algorithm. *IEEE Access, 9*, 124044-124058.

[15] Hao, J., Zhu, C., & Guo, X. (2021). Wind Power Short-Term Forecasting Model Based on the Hierarchical Output Power and Poisson Re-Sampling Random Forest Algorithm. *IEEE Access, 9*, 6478-6487.

[16] Govindarajan, S., Ardila-Rey, J.A., Krithivasan, K., Subbaiah, J., Sannidhi, N., & Balasubramanian, M. (2021). Development of Hypergraph Based Improved Random Forest Algorithm for Partial Discharge Pattern Classification. *IEEE Access, 9*, 96-109.

[17] Jayasinghe, W.J., Deo, R.C., Ghahramani, A., Ghimire, S., & Raj, N. (2021). Deep multi-stage reference evapotranspiration forecasting model: Multivariate empirical mode decomposition integrated with Boruta-random forest algorithm. *IEEE Access, PP*, 1-1.

[18] Zhu, M., Xia, J., Jin, X., Yan, M., Cai, G., Yan, J., & Ning, G. (2018). Class Weights Random Forest Algorithm for Processing Class Imbalanced Medical Data. *IEEE Access, 6*, 4641-4652.

[19] Jia, D.Y., Li, Z., & Zhang, C. (2020). A Parametric Optimization Oriented, AFSA Based Random Forest Algorithm: Application to the Detection of Cervical Epithelial Cells. *IEEE Access, 8*, 64891-64905.

[20] Li, X., Gao, Y., Zhang, H., & Liao, Y. (2020). Passenger Travel Behavior in Public Transport Corridor After the Operation of Urban Rail Transit: A Random Forest Algorithm Approach. *IEEE Access, 8*, 211303-211314.