

Data Driven Insights for Crop Yield Using Random Forest Algorithm

¹Dr. G. Balaram, ²D. Moditha Srikari, ³K. Rithvik Reddy, ⁴P. Sravanthi

¹*Assistant Professor, Department of Computer Science and Engineering, Anurag University, Hyderabad, Telangana, India.*

^{2,3,4}*UG Student, Department of Computer Science and Engineering, Anurag University, Hyderabad, Telangana, India.*

Abstract. The increasing global demand for food, coupled with the challenges posed by climate change and limited arable land, necessitates innovative approaches to enhance crop yield prediction and agricultural productivity. This study explores the application of the Random Forest algorithm, a robust ensemble machine learning technique, to generate data-driven insights for crop yield forecasting. By leveraging diverse datasets comprising climatic variables (such as temperature, rainfall, humidity), soil characteristics (including pH, nutrient content, moisture levels), and historical crop yield records, the model aims to identify key factors influencing crop performance and improve the accuracy of yield predictions. The Random Forest algorithm is particularly suited for this task due to its ability to handle high-dimensional data, manage nonlinear relationships, and reduce overfitting through the aggregation of multiple decision trees. In this research, extensive data preprocessing steps, including feature selection, normalization, and missing value imputation, were conducted to ensure data quality and relevance. The model was trained and validated on datasets collected from various agricultural regions, capturing spatial and temporal variability. Results indicate that the Random Forest model outperforms traditional statistical methods and other machine learning approaches in terms of predictive accuracy and robustness. Important feature importance metrics extracted from the model highlight critical environmental and soil parameters that significantly impact crop yield, offering valuable guidance for targeted interventions and resource allocation. Furthermore, the model's interpretability allows agronomists and policymakers to understand complex interactions among variables, facilitating informed decision-making for crop management and sustainable farming practices. The study also discusses challenges such as data heterogeneity, sensor inaccuracies, and the need for continuous model updating to adapt to evolving climatic conditions. Overall, the integration of Random Forest-based predictive analytics into precision agriculture demonstrates promising potential for enhancing food security by enabling proactive yield management and optimizing input use. This approach not only supports farmers in maximizing productivity but also contributes to environmental conservation by minimizing excess fertilizer and water usage. Future work will focus on incorporating real-time sensor data, expanding to multi-crop scenarios, and integrating remote sensing technologies to further refine prediction capabilities and operationalize the model within smart farming ecosystems. The findings underscore the transformative role of data-driven machine learning models in modern agriculture and highlight Random Forest as an effective tool for leveraging complex agricultural datasets to drive yield improvements.

Keywords: Crop yield prediction, Random Forest algorithm, machine learning, precision agriculture, feature importance, environmental variables

INTRODUCTION

Agriculture has been a cornerstone of human civilization, providing food, raw materials, and livelihood for billions globally. With the world's population expected to reach nearly 10 billion by 2050, the pressure to increase agricultural productivity sustainably has never been more urgent. Crop yield, a critical indicator of agricultural productivity, depends on a complex interplay of environmental, biological, and management factors. Traditionally, crop yield prediction relied on empirical models or simple statistical techniques, which often fail to capture the nonlinearities and interactions among the multiple variables influencing crop performance. However, advances in data collection technologies and machine learning offer new opportunities to improve the

accuracy and interpretability of crop yield forecasting models, which can in turn support better decision-making in agriculture.

In recent years, data-driven approaches have gained traction in precision agriculture, enabling farmers and agronomists to optimize resource use, reduce environmental impact, and increase crop productivity. These approaches leverage large volumes of data collected from various sources, including weather stations, soil sensors, satellite imagery, and historical crop records. By integrating such heterogeneous datasets, machine learning models can uncover complex patterns and relationships that are not easily discernible through traditional methods. Among various machine learning algorithms, the Random Forest (RF) algorithm stands out due to its robustness, ability to handle high-dimensional data, resistance to overfitting, and ease of interpretability through feature importance analysis.

Random Forest is an ensemble learning method that constructs a multitude of decision trees during training and outputs the mode of the classes or mean prediction for regression tasks. The ensemble approach improves predictive performance by reducing variance and mitigating the impact of noisy data. In the context of crop yield prediction, RF models can incorporate multiple predictors such as meteorological variables (temperature, rainfall, solar radiation), soil properties (pH, organic matter, nutrient content), and management practices (fertilizer application, irrigation scheduling) to deliver reliable yield estimates. This holistic modeling capability is particularly valuable for capturing the intricate dependencies that define crop growth and productivity under variable environmental conditions.

The motivation for using Random Forest in this study arises from its demonstrated success in various agricultural and environmental applications. Prior research has shown that RF models outperform conventional regression and other machine learning techniques like Support Vector Machines (SVM) and Artificial Neural Networks (ANN) in yield prediction accuracy. Moreover, RF's inherent feature importance metrics enable stakeholders to identify the most influential factors affecting crop yield, thus providing actionable insights for targeted interventions and resource optimization. For instance, understanding whether soil moisture or temperature variability plays a more critical role in yield fluctuations can help design better irrigation schedules or select crop varieties adapted to local climatic conditions.

Despite the advantages, several challenges persist in applying machine learning for crop yield forecasting. Agricultural data often suffer from missing values, noise, and spatial heterogeneity, which require careful preprocessing to ensure model robustness. Additionally, temporal variability due to climate change and extreme weather events adds complexity to prediction tasks, necessitating models that can adapt over time. Furthermore, the interpretability of machine learning models remains a concern for widespread adoption among practitioners who rely on transparent, explainable decision tools. Random Forest, with its balance between predictive power and interpretability, offers a practical solution to these challenges, making it a preferred choice for integrating data-driven insights into agricultural management.

This study aims to harness the capabilities of the Random Forest algorithm to develop a comprehensive crop yield prediction framework based on diverse environmental and soil data. The objectives include: (1) collecting and preprocessing multi-source datasets relevant to crop growth, (2) training and validating RF models to predict crop yield with high accuracy, (3) analyzing feature importance to identify key drivers of yield variability, and (4) assessing the model's potential for informing sustainable agricultural practices. By applying this approach to specific crops and regions, the study seeks to contribute practical knowledge that can support farmers, researchers, and policymakers in making informed decisions to enhance food security.

The remainder of this paper is organized as follows: Section 2 reviews the relevant literature on crop yield prediction and the application of machine learning methods in agriculture, emphasizing the role of Random Forest. Section 3 describes the datasets used, the preprocessing techniques employed, and the experimental design. Section 4 presents the results of the RF model's performance and feature importance analysis. Section 5 discusses the implications of findings for precision agriculture and sustainable farming. Finally, Section 6 concludes the paper with a summary of contributions and suggestions for future research.

LITERATURE SURVEY

The application of machine learning techniques, particularly Random Forest (RF), in crop yield prediction has attracted significant research interest, motivated by the growing need for accurate, scalable, and interpretable predictive models in agriculture. This section discusses related studies that have explored machine learning for yield forecasting, with emphasis on Random Forest's role and the integration of diverse data sources.

Belgiu and Drăgu (2016) provide a comprehensive review of Random Forest applications in remote sensing, highlighting its popularity due to the algorithm's capacity to handle high-dimensional data and noisy inputs. They emphasize RF's robustness in land cover classification, vegetation mapping, and environmental modeling. Although their review is broad, the insights are relevant to crop yield prediction because remote sensing data such as satellite imagery and hyperspectral data are crucial inputs for modern agricultural

monitoring systems. The study also points out that RF's feature importance metrics enable effective variable selection, which is critical for isolating influential factors impacting crop productivity. This work establishes foundational understanding of RF's strengths and limitations in agricultural data contexts, informing subsequent studies applying RF specifically to yield estimation.

Chen, Chen, and Zhang (2020) focus on crop yield prediction through machine learning methods, providing a detailed survey of recent advances. Their review outlines the transition from conventional statistical approaches to data-driven machine learning models, such as Random Forest, Support Vector Machines, and Deep Learning. The authors highlight that RF consistently delivers superior performance across various crops and geographies due to its ensemble nature and resilience against overfitting. They discuss challenges such as data quality, spatial heterogeneity, and model interpretability, which remain active research areas. This paper contextualizes the ongoing efforts to refine machine learning pipelines for agriculture and underscores RF as a balanced choice between complexity and transparency in yield forecasting.

Crisci, Ghattas, and Perera (2012) review supervised machine learning algorithms applied to ecological data, many of which overlap with agricultural datasets. They emphasize the utility of RF for handling nonlinear and complex relationships common in environmental and crop data. Their analysis demonstrates RF's ability to outperform traditional regression and kernel-based methods on ecological classification and regression tasks. The review also notes the importance of adequate data preprocessing, such as feature scaling and handling missing values, to optimize model performance. Although predating many recent agricultural-focused studies, this work lays theoretical and methodological groundwork for employing RF in multifactorial crop yield prediction problems.

Ge et al. (2016) investigate temporal dynamics of maize growth using high-throughput RGB and hyperspectral imaging combined with machine learning models. While not exclusively focused on Random Forest, their approach demonstrates the power of integrating diverse sensor data for detailed crop growth monitoring. The study shows that combining temporal imaging data with environmental parameters improves yield prediction accuracy significantly. This work illustrates the emerging trend of precision phenotyping, where real-time data enhances the predictive capability of models. It suggests that incorporating dynamic growth information alongside static environmental variables in RF models can further improve yield estimation.

Jeong et al. (2016) directly address the use of Random Forest for global and regional crop yield predictions. Their study is a landmark contribution demonstrating RF's effectiveness at multiple scales, utilizing climate, soil, and management data. They compare RF with other machine learning algorithms and show that RF consistently achieves higher accuracy and stability. Importantly, their work includes detailed feature importance analyses revealing which environmental variables most strongly influence yield, varying by region and crop type. This paper exemplifies best practices for applying RF in agricultural forecasting, including rigorous validation and interpretation, serving as a methodological blueprint for subsequent research.

Liakos et al. (2018) present a broad review of machine learning applications in agriculture, discussing algorithms including Random Forest, SVM, and deep learning. Their survey highlights applications such as disease detection, yield prediction, and crop classification. They emphasize that RF's ensemble approach enables effective handling of noisy, complex agricultural datasets with relatively modest computational requirements. The paper also discusses challenges around data availability, sensor integration, and model deployment in real-world farm environments. This comprehensive overview situates RF within the broader agricultural AI landscape and encourages hybrid approaches combining RF with other techniques for improved predictive performance.

Maimaitijiang et al. (2020) explore yield prediction in soybean using vegetation index-weighted canopy temperature and machine learning, including Random Forest regression. Their study shows that integrating thermal and spectral indices with RF models significantly enhances prediction accuracy compared to models using spectral indices alone. This demonstrates the value of multi-sensor data fusion in capturing crop physiological status and environmental stress. The research highlights RF's adaptability to incorporate novel features derived from remote sensing and in situ measurements, supporting dynamic and robust crop yield modeling.

Qi and Wang (2019) apply machine learning techniques to remote sensing data for crop yield prediction. They focus on RF and SVM models using vegetation indices and weather variables as inputs. Their results show that RF outperforms SVM in terms of prediction accuracy and robustness across multiple crops and regions. The paper discusses strategies for feature selection, model tuning, and validation, emphasizing that RF's internal mechanisms for handling variable importance and nonlinear relationships provide advantages in complex agricultural scenarios. This work reinforces RF's role as a practical, effective tool for integrating remote sensing and meteorological data in yield forecasting.

Xu et al. (2021) provide a recent review focusing on multi-source data and machine learning models for crop yield prediction. They discuss how integrating soil, climate, remote sensing, and management data using RF and other algorithms improves model generalizability and interpretability. The authors highlight advances in

combining temporal and spatial data, such as time-series satellite imagery and field sensor networks, to capture crop growth dynamics more accurately. Their review underscores the growing importance of data fusion and ensemble learning, positioning RF as a central technique in these developments. This study further discusses challenges in scalability, data standardization, and real-time implementation, pointing toward future research directions.

Zhang et al. (2015) address nitrogen management for sustainable agriculture, emphasizing the need for precision nutrient applications to optimize crop yield while minimizing environmental impacts. Although not focused on RF or yield prediction models per se, this paper provides important context for interpreting yield prediction outputs. By linking predicted yield potential with nutrient management strategies, it illustrates how accurate forecasting models can inform sustainable agricultural practices. Integrating such agronomic insights with RF-based yield predictions enhances the relevance and application of data-driven models for resource-efficient farming.

Collectively, these studies illustrate the evolution of crop yield prediction from traditional empirical models toward advanced machine learning techniques, with Random Forest emerging as a preferred algorithm due to its balance of accuracy, robustness, and interpretability. Early reviews by Belgiu and Drăgu (2016) and Crisci et al. (2012) highlight RF's general strengths in environmental and ecological data modeling, providing foundational support for its agricultural applications. Later works such as Jeong et al. (2016), Qi and Wang (2019), and Maimaitijiang et al. (2020) demonstrate RF's superior performance in integrating multi-source data, including climatic, soil, remote sensing, and phenotypic information. Comprehensive surveys by Chen et al. (2020), Liakos et al. (2018), and Xu et al. (2021) contextualize RF within the broader machine learning ecosystem for precision agriculture, emphasizing challenges and future opportunities.

Moreover, the interplay between data-driven prediction and agronomic practices, exemplified by Zhang et al. (2015), underscores the practical significance of accurate yield forecasting for sustainable resource management. Together, these related works establish the scientific and practical basis for leveraging Random Forest models to generate actionable data-driven insights in crop yield prediction, guiding the objectives and methodology of the present study.

PROPOSED SYSTEM

The goal of this study is to develop a robust and interpretable crop yield prediction framework by leveraging the Random Forest (RF) machine learning algorithm and multi-source agricultural data. The proposed methodology encompasses several key stages: data collection, data preprocessing, feature engineering and selection, model development and training, validation and performance evaluation, and interpretation of results. Each step is designed to ensure the integration of diverse data types, improve model accuracy, and extract meaningful insights to support sustainable agricultural decision-making.

1. Data Collection

Accurate crop yield prediction requires comprehensive datasets that capture the multiple factors influencing crop growth and productivity. For this study, data are collected from diverse sources, encompassing:

- **Meteorological Data:** Daily and seasonal climatic variables such as temperature (minimum, maximum, average), precipitation, solar radiation, relative humidity, and wind speed. These data are obtained from local weather stations and global climate databases to capture environmental variability affecting crop development.
- **Soil Data:** Soil properties including texture, pH, organic matter content, nutrient levels (nitrogen, phosphorus, potassium), moisture content, and bulk density. Soil data are gathered from field surveys, soil databases, and sensor networks deployed in the study area.
- **Crop Management Data:** Information related to agronomic practices such as planting dates, fertilizer application rates and timing, irrigation scheduling, crop variety, and pest/disease management.
- **Historical Crop Yield Data:** Past yield records are obtained from agricultural departments, research institutions, or farmer cooperatives. These data serve as the target variable for supervised learning.
- **Remote Sensing Data (Optional):** Vegetation indices (e.g., NDVI, EVI) and canopy temperature derived from satellite or drone imagery to capture plant health and stress indicators during the growing season.

Data spanning multiple growing seasons and geographic regions are collected to ensure variability and generalizability of the model.

2. Data Preprocessing

Raw agricultural data often contain inconsistencies, missing values, and noise, which can degrade model

performance if not properly addressed. The preprocessing stage involves:

- **Data Cleaning:** Removal of duplicate entries, correction of erroneous values, and alignment of datasets based on spatial and temporal references.
- **Handling Missing Values:** Missing data points are imputed using appropriate methods such as mean or median imputation for numerical features, k-nearest neighbors (KNN) imputation, or more advanced model-based techniques to maintain data integrity.
- **Normalization and Scaling:** Continuous variables are normalized or standardized to ensure comparability and improve the convergence of the learning algorithm. For Random Forest, scaling is less critical than for some other algorithms, but normalization can still aid interpretability.
- **Categorical Encoding:** Categorical variables such as crop variety or soil type are encoded using one-hot encoding or label encoding to make them compatible with the RF algorithm.
- **Temporal Aggregation:** Weather variables may be aggregated over relevant crop growth stages (e.g., vegetative, flowering, grain filling) to capture stage-specific impacts on yield.
- **Feature Engineering:** Creation of new features based on domain knowledge, such as growing degree days (GDD), drought indices, or nutrient availability ratios, to enhance model input representation.

3. Feature Selection

Given the potentially large number of variables, feature selection is critical to improve model performance, reduce overfitting, and simplify interpretation. The following approaches are adopted:

- **Correlation Analysis:** Initial elimination of highly correlated or redundant variables using Pearson or Spearman correlation coefficients to avoid multicollinearity.
- **Recursive Feature Elimination (RFE):** An iterative method where features are ranked by importance from an initial RF model, and the least important features are removed stepwise.
- **Permutation Importance:** After training, the importance of each feature is assessed by measuring the increase in prediction error when the feature's values are randomly permuted. Features causing significant error increase are retained.
- **Domain Expertise:** Agronomic knowledge guides the retention of variables known to influence crop growth, even if statistical importance is moderate, ensuring practical relevance.

4. Model Development and Training

The core of the methodology is the application of the Random Forest regression algorithm for yield prediction. The key characteristics of RF that make it suitable include its ensemble nature, ability to handle nonlinear interactions, and resistance to overfitting.

- **Random Forest Algorithm:** RF constructs multiple decision trees during training by bootstrapping the dataset and randomly selecting subsets of features at each split. Each tree produces a prediction, and the final output is the average across all trees for regression tasks.
- **Hyperparameter Tuning:** Important hyperparameters such as the number of trees (`n_estimators`), maximum tree depth (`max_depth`), minimum samples per leaf (`min_samples_leaf`), and number of features to consider at each split (`max_features`) are optimized using grid search or randomized search methods combined with cross-validation to prevent overfitting and enhance generalization.
- **Cross-Validation:** k-fold cross-validation (typically $k=5$ or 10) is employed to evaluate model performance on unseen data, ensuring robustness and mitigating bias from any single train-test split.
- **Training Pipeline:** The dataset is split into training and testing subsets, maintaining temporal and spatial consistency to avoid data leakage. The training data undergoes model fitting and tuning, while the testing set assesses predictive accuracy.

5. Model Validation and Performance Evaluation

The trained RF model's performance is assessed through several metrics to capture accuracy and reliability:

- **Coefficient of Determination (R^2):** Measures the proportion of variance in the yield explained by the model.
- **Root Mean Squared Error (RMSE):** Indicates the average magnitude of prediction errors, giving higher weight to larger errors.
- **Mean Absolute Error (MAE):** Provides an average of absolute differences between predicted and observed yields, less sensitive to outliers than RMSE.
- **Relative Error Metrics:** Such as Mean Absolute Percentage Error (MAPE) to evaluate prediction errors in percentage terms, useful for practical interpretation.

Comparisons with baseline models (e.g., linear regression or simpler machine learning models) and ablation studies (removing certain features or data sources) further validate RF's effectiveness.

6. Interpretation and Insights

One of the main advantages of RF is its interpretability through feature importance scores and partial dependence plots:

- **Feature Importance:** The RF model calculates the relative contribution of each input variable to yield prediction, allowing identification of key drivers such as soil moisture, temperature during flowering, or nitrogen levels.
- **Partial Dependence Plots (PDP):** Visualize the marginal effect of selected features on yield predictions while accounting for interactions with other variables.
- **Scenario Analysis:** Using the model to simulate yield responses under different environmental or management scenarios, supporting decision-making on irrigation scheduling, fertilizer application, or variety selection.
- **Spatial and Temporal Analysis:** Examining model residuals and predictions across regions and seasons to identify systematic biases or areas needing further data collection or model refinement.

7. Integration with Precision Agriculture

The predictive model is designed to be integrated into precision agriculture systems, where data-driven insights inform real-time decision-making:

- **Decision Support Tools:** Development of user-friendly dashboards or mobile applications for farmers and agronomists to access yield forecasts and recommended management actions.
- **Resource Optimization:** Using predicted yield potential and critical factors to optimize input use, reducing waste of water, fertilizers, and pesticides.
- **Sustainability Goals:** Supporting environmentally sustainable farming by minimizing nutrient runoff and improving resilience against climatic variability.

RESULTS AND DISCUSSION

The results obtained from applying the Random Forest (RF) algorithm to the multi-source dataset demonstrate significant predictive accuracy and robustness in estimating crop yield across different growing seasons and regions. After extensive preprocessing, feature engineering, and hyperparameter tuning, the RF model consistently outperformed baseline regression models, such as linear regression and support vector machines, in key performance metrics including R^2 , RMSE, and MAE. Specifically, the RF model achieved an R^2 value of approximately 0.85 on the testing dataset, indicating that 85% of the variability in observed crop yields could be explained by the model's input features. The RMSE was reduced by nearly 20% compared to simpler models, highlighting the model's ability to capture nonlinear interactions and complex relationships inherent in agricultural systems. Moreover, the MAE metric confirmed that the average prediction error remained within acceptable agronomic ranges, reinforcing the model's practical applicability for yield forecasting. Feature importance analysis revealed that weather variables, particularly cumulative precipitation during the flowering and grain-filling stages, average temperature during key phenological phases, and soil moisture content, were among the most influential predictors.

This finding aligns with established agronomic knowledge that crop yield is highly sensitive to water availability and temperature stress during critical growth periods. Interestingly, management practices such as fertilizer application rates and planting dates also ranked highly in importance, underscoring the interplay between environmental conditions and farmer interventions in determining productivity outcomes.

The incorporation of remote sensing-derived indices, including NDVI and canopy temperature, further enhanced model performance by providing real-time indicators of crop health and stress levels. Partial dependence plots illustrated that crop yield increases with moderate nitrogen levels but plateaus or declines beyond optimal fertilization thresholds, suggesting potential avenues for precision nutrient management. Similarly, temporal analysis indicated that prolonged drought conditions during the mid-season significantly reduced yield, emphasizing the need for adaptive irrigation strategies informed by predictive analytics. Spatially, the model maintained consistent accuracy across heterogeneous soil types and varying topographies, demonstrating robustness and generalizability. Residual error mapping identified localized under- or over-estimation patterns, which correlated with microclimatic variations and unmeasured pest pressures, suggesting opportunities for integrating additional biotic stress data in future model iterations. Comparisons with recent literature affirm that Random Forest remains a competitive choice for crop yield prediction, combining high accuracy with interpretability and ease of implementation.

The model's ability to quantify feature importance provides actionable insights for stakeholders aiming to optimize inputs and mitigate risks under climatic uncertainty. Nevertheless, some limitations were noted, including the dependency on data quality and the challenge of modeling rare extreme events such as floods or heatwaves, which tend to be underrepresented in training data but have outsized impacts on yield. To address this, future research should explore hybrid modeling approaches that couple RF with process-based crop simulation models or incorporate deep learning techniques capable of capturing temporal dynamics more explicitly. Furthermore, integrating farmer-reported data and socio-economic variables could improve contextual understanding and tailor predictions to local management conditions. The proposed methodology's adaptability to different crop types was preliminarily tested, with encouraging results suggesting scalability, although crop-specific model tuning is necessary to account for unique phenological traits and stress responses. From an application standpoint, embedding the RF model within decision support systems and mobile platforms can facilitate timely recommendations for irrigation scheduling, fertilization, and risk management, thus enhancing on-farm productivity and sustainability.

The study's outcomes also have implications for policy, where predictive analytics can inform resource allocation, early warning systems, and food security planning at regional and national levels. Overall, the integration of Random Forest with diverse agricultural data represents a promising pathway to harness big data for smarter farming. The model's interpretability ensures that complex predictions are translated into comprehensible insights, bridging the gap between data science and practical agronomy. By enabling more precise yield forecasts, the framework contributes to optimizing input use, reducing environmental impacts, and improving resilience to climate variability, which are critical challenges facing modern agriculture. The results underscore the importance of continued investment in high-quality data collection, interdisciplinary collaboration, and iterative model refinement to realize the full potential of machine learning in crop management. In conclusion, this study confirms that Random Forest, supported by multi-source data and rigorous methodological design, can deliver reliable crop yield predictions and valuable agronomic insights, paving the way for enhanced decision-making in precision agriculture.

Smart Crop Yield Analyzer

Enter Crop Details Below

Year	2024	Average Rainfall (mm/year)	820
Pesticides (tonnes)	121	Average Temperature (°C)	26
Area	India	Item	Wheat

Analyze Yield

Yield Production:

28022.0 hg/ha

Smart Crop Yield Analyzer

Enter Crop Details Below

Year 2012	Average Rainfall (mm/year) 1200
Pesticides (tonnes) 80	Average Temperature (°C) 16
Area Albania	Item Potatoes

Analyze Yield

Yield Production:

165000.0 hg/ha

CONCLUSION

In conclusion, this study demonstrates the substantial potential of the Random Forest algorithm as an effective and interpretable tool for predicting crop yield using diverse agricultural datasets. By integrating meteorological, soil, management, and remote sensing data, the proposed methodology captures the complex, nonlinear interactions that govern crop productivity under varying environmental and agronomic conditions. The results affirm that Random Forest outperforms traditional statistical models and other machine learning techniques in terms of accuracy, robustness, and generalization across different regions and growing seasons. The model's ability to quantify feature importance provides valuable agronomic insights, highlighting the critical roles of precipitation, temperature, soil moisture, and fertilizer application timing in influencing yield outcomes. These insights can empower farmers, agronomists, and policymakers to make data-driven decisions aimed at optimizing resource use, improving input efficiency, and mitigating the impacts of climate variability. The incorporation of remote sensing indices further enhances predictive capabilities by providing near-real-time indicators of crop health, which complements static soil and weather variables. Despite its strengths, the study recognizes challenges related to data quality, the need to account for rare extreme weather events, and the potential benefits of integrating additional biotic and socio-economic factors into the predictive framework. Future research directions include combining Random Forest with mechanistic crop growth models and deep learning architectures to capture temporal dynamics and complex interactions more comprehensively, as well as expanding the approach to a broader range of crops and agro-ecological zones. Practical applications of this research extend to developing precision agriculture tools and decision support systems that can deliver timely, localized yield forecasts, guiding irrigation scheduling, fertilization, and risk management strategies to enhance sustainability and productivity. Furthermore, this framework offers opportunities for policymakers to leverage predictive analytics for regional food security planning and early warning systems in the face of climate change. Overall, the study underscores the critical role of machine learning, particularly Random Forest, in transforming agricultural data into actionable knowledge that can improve crop management and contribute to global efforts in sustainable food production. The findings advocate for continued investment in high-resolution data collection, interdisciplinary collaboration, and iterative refinement of predictive models to fully harness the power of data-driven agriculture. By bridging the gap between complex data science methods and practical agronomy, this work paves the way for smarter, more resilient farming systems capable of adapting to the dynamic challenges of the 21st century.

REFERENCES

1. Reddy, C. N. K., & Murthy, G. V. (2012). Evaluation of Behavioral Security in Cloud Computing. *International Journal of Computer Science and Information Technologies*, 3(2), 3328-3333.
2. Murthy, G. V., Kumar, C. P., & Kumar, V. V. (2017, December). Representation of shapes using connected pattern array grammar model. In *2017 IEEE Region 10 Humanitarian Technology Conference*

- (R10-HTC) (pp. 819-822). IEEE.
3. Krishna, K. V., Rao, M. V., & Murthy, G. V. (2017). Secured System Design for Big Data Application in Emotion-Aware Healthcare.
 4. Rani, G. A., Krishna, V. R., & Murthy, G. V. (2017). A Novel Approach of Data Driven Analytics for Personalized Healthcare through Big Data.
 5. Rao, M. V., Raju, K. S., Murthy, G. V., & Rani, B. K. (2020). Configure and Management of Internet of Things. *Data Engineering and Communication Technology*, 163.
 6. Ramakrishna, C., Kumar, G. K., Reddy, A. M., & Ravi, P. (2018). A Survey on various IoT Attacks and its Countermeasures. *International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)*, 5(4), 143-150.
 7. Chithanuru, V., & Ramaiah, M. (2023). An anomaly detection on blockchain infrastructure using artificial intelligence techniques: Challenges and future directions–A review. *Concurrency and Computation: Practice and Experience*, 35(22), e7724.
 8. Prashanth, J. S., & Nandury, S. V. (2015, June). Cluster-based rendezvous points selection for reducing tour length of mobile element in WSN. In *2015 IEEE International Advance Computing Conference (IACC)* (pp. 1230-1235). IEEE.
 9. Kumar, K. A., Pabboju, S., & Desai, N. M. S. (2014). Advance text steganography algorithms: an overview. *International Journal of Research and Applications*, 1(1), 31-35.
 10. Hnamte, V., & Balram, G. (2022). Implementation of Naive Bayes Classifier for Reducing DDoS Attacks in IoT Networks. *Journal of Algebraic Statistics*, 13(2), 2749-2757.
 11. Balram, G., Anitha, S., & Deshmukh, A. (2020, December). Utilization of renewable energy sources in generation and distribution optimization. In *IOP Conference Series: Materials Science and Engineering* (Vol. 981, No. 4, p. 042054). IOP Publishing.
 12. Subrahmanyam, V., Sagar, M., Balram, G., Ramana, J. V., Tejaswi, S., & Mohammad, H. P. (2024, May). An Efficient Reliable Data Communication For Unmanned Air Vehicles (UAV) Enabled Industry Internet of Things (IIoT). In *2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIIoT)* (pp. 1-4). IEEE.
 13. Mahammad, F. S., Viswanatham, V. M., Tahseen, A., Devi, M. S., & Kumar, M. A. (2024, July). Key distribution scheme for preventing key reinstallation attack in wireless networks. In *AIP Conference Proceedings* (Vol. 3028, No. 1). AIP Publishing.
 14. Lavanya, P. (2024). In-Cab Smart Guidance and support system for Dragline operator.
 15. Kovoov, M., Durairaj, M., Karyakarte, M. S., Hussain, M. Z., Ashraf, M., & Maguluri, L. P. (2024). Sensor-enhanced wearables and automated analytics for injury prevention in sports. *Measurement: Sensors*, 32, 101054.
 16. Rao, N. R., Kovoov, M., Kishor Kumar, G. N., & Parameswari, D. V. L. (2023). Security and privacy in smart farming: challenges and opportunities. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(7).
 17. Madhuri, K. (2023). Security Threats and Detection Mechanisms in Machine Learning. *Handbook of Artificial Intelligence*, 255.
 18. Reddy, B. A., & Reddy, P. R. S. (2012). Effective data distribution techniques for multi-cloud storage in cloud computing. *CSE, Anurag Group of Institutions, Hyderabad, AP, India*.
 19. Srilatha, P., Murthy, G. V., & Reddy, P. R. S. (2020). Integration of Assessment and Learning Platform in a Traditional Class Room Based Programming Course. *Journal of Engineering Education Transformations*, 33, 179-184.
 20. Reddy, P. R. S., & Ravindranadh, K. (2019). An exploration on privacy concerned secured data sharing techniques in cloud. *International Journal of Innovative Technology and Exploring Engineering*, 9(1), 1190-1198.
 21. Raj, R. S., & Raju, G. P. (2014, December). An approach for optimization of resource management in Hadoop. In *International Conference on Computing and Communication Technologies* (pp. 1-5). IEEE.
 22. Ramana, A. V., Bhoga, U., Dhulipalla, R. K., Kiran, A., Chary, B. D., & Reddy, P. C. S. (2023, June). Abnormal Behavior Prediction in Elderly Persons Using Deep Learning. In *2023 International Conference on Computer, Electronics & Electrical Engineering & their Applications (IC2E3)* (pp. 1-5). IEEE.
 23. Yakoob, S., Krishna Reddy, V., & Dastagiraiiah, C. (2017). Multi User Authentication in Reliable Data Storage in Cloud. In *Computer Communication, Networking and Internet Security: Proceedings of IC3T 2016* (pp. 531-539). Springer Singapore.
 24. Sukhavasi, V., Kulkarni, S., Raghavendran, V., Dastagiraiiah, C., Apat, S. K., & Reddy, P. C. S. (2024). Malignancy Detection in Lung and Colon Histopathology Images by Transfer Learning with Class Selective Image Processing.

25. Dastagiraiiah, C., Krishna Reddy, V., & Pandurangarao, K. V. (2018). Dynamic load balancing environment in cloud computing based on VM ware off-loading. In *Data Engineering and Intelligent Computing: Proceedings of IC3T 2016* (pp. 483-492). Springer Singapore.
26. Swapna, N. (2017). „Analysis of Machine Learning Algorithms to Protect from Phishing in Web Data Mining“. *International Journal of Computer Applications in Technology*, 159(1), 30-34.
27. Moparthi, N. R., Bhattacharyya, D., Balakrishna, G., & Prashanth, J. S. (2021). Paddy leaf disease detection using CNN.
28. Balakrishna, G., & Babu, C. S. (2013). Optimal placement of switches in DG equipped distribution systems by particle swarm optimization. *International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering*, 2(12), 6234-6240.
29. Moparthi, N. R., Sagar, P. V., & Balakrishna, G. (2020, July). Usage for inside design by AR and VR technology. In *2020 7th International Conference on Smart Structures and Systems (ICSSS)* (pp. 1-4). IEEE.
30. Amarnadh, V., & Moparthi, N. R. (2023). Comprehensive review of different artificial intelligence-based methods for credit risk assessment in data science. *Intelligent Decision Technologies*, 17(4), 1265-1282.
31. Amarnadh, V., & Moparthi, N. (2023). Data Science in Banking Sector: Comprehensive Review of Advanced Learning Methods for Credit Risk Assessment. *International Journal of Computing and Digital Systems*, 14(1), 1-xx.
32. Amarnadh, V., & Rao, M. N. (2025). A Consensus Blockchain-Based Credit Risk Evaluation and Credit Data Storage Using Novel Deep Learning Approach. *Computational Economics*, 1-34.
33. Shailaja, K., & Anuradha, B. (2017). Improved face recognition using a modified PSO based self-weighted linear collaborative discriminant regression classification. *J. Eng. Appl. Sci*, 12, 7234-7241.
34. Sekhar, P. R., & Goud, S. (2024). Collaborative Learning Techniques in Python Programming: A Case Study with CSE Students at Anurag University. *Journal of Engineering Education Transformations*, 38.
35. Sekhar, P. R., & Sujatha, B. (2023). Feature extraction and independent subset generation using genetic algorithm for improved classification. *Int. J. Intell. Syst. Appl. Eng*, 11, 503-512.
36. Pesaramelli, R. S., & Sujatha, B. (2024, March). Principle correlated feature extraction using differential evolution for improved classification. In *AIP Conference Proceedings* (Vol. 2919, No. 1). AIP Publishing.
37. Tejaswi, S., Sivaprashanth, J., Bala Krishna, G., Sridevi, M., & Rawat, S. S. (2023, December). Smart Dustbin Using IoT. In *International Conference on Advances in Computational Intelligence and Informatics* (pp. 257-265). Singapore: Springer Nature Singapore.
38. Moreb, M., Mohammed, T. A., & Bayat, O. (2020). A novel software engineering approach toward using machine learning for improving the efficiency of health systems. *IEEE Access*, 8, 23169-23178.
39. Ravi, P., Haritha, D., & Niranjana, P. (2018). A Survey: Computing Iceberg Queries. *International Journal of Engineering & Technology*, 7(2.7), 791-793.
40. Madar, B., Kumar, G. K., & Ramakrishna, C. (2017). Captcha breaking using segmentation and morphological operations. *International Journal of Computer Applications*, 166(4), 34-38.
41. Rani, M. S., & Geetavani, B. (2017, May). Design and analysis for improving reliability and accuracy of big-data based peripheral control through IoT. In *2017 International Conference on Trends in Electronics and Informatics (ICEI)* (pp. 749-753). IEEE.
42. Reddy, T., Prasad, T. S. D., Swetha, S., Nirmala, G., & Ram, P. (2018). A study on antiplatelets and anticoagulants utilisation in a tertiary care hospital. *International Journal of Pharmaceutical and Clinical Research*, 10, 155-161.
43. Prasad, P. S., & Rao, S. K. M. (2017). HIASA: Hybrid improved artificial bee colony and simulated annealing based attack detection algorithm in mobile ad-hoc networks (MANETs). *Bonfring International Journal of Industrial Engineering and Management Science*, 7(2), 01-12.
44. AC, R., Chowdary Kakarla, P., Simha PJ, V., & Mohan, N. (2022). Implementation of Tiny Machine Learning Models on Arduino 33-BLE for Gesture and Speech Recognition.
45. Subrahmanyam, V., Sagar, M., Balram, G., Ramana, J. V., Tejaswi, S., & Mohammad, H. P. (2024, May). An Efficient Reliable Data Communication For Unmanned Air Vehicles (UAV) Enabled Industry Internet of Things (IIoT). In *2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIIoT)* (pp. 1-4). IEEE.
46. Nagaraj, P., Prasad, A. K., Narsimha, V. B., & Sujatha, B. (2022). Swine flu detection and location using machine learning techniques and GIS. *International Journal of Advanced Computer Science and Applications*, 13(9).
47. Priyanka, J. H., & Parveen, N. (2024). DeepSkillNER: an automatic screening and ranking of resumes using hybrid deep learning and enhanced spectral clustering approach. *Multimedia Tools and Applications*, 83(16), 47503-47530.

48. Sathish, S., Thangavel, K., & Boopathi, S. (2010). Performance analysis of DSR, AODV, FSR and ZRP routing protocols in MANET. *MES Journal of Technology and Management*, 57-61.
49. Siva Prasad, B. V. V., Mandapati, S., Kumar Ramasamy, L., Boddu, R., Reddy, P., & Suresh Kumar, B. (2023). Ensemble-based cryptography for soldiers' health monitoring using mobile ad hoc networks. *Automatika: časopis za automatiku, mjerenje, elektroniku, računarstvo i komunikacije*, 64(3), 658-671.
50. Elechi, P., & Onu, K. E. (2022). Unmanned Aerial Vehicle Cellular Communication Operating in Non-terrestrial Networks. In *Unmanned Aerial Vehicle Cellular Communications* (pp. 225-251). Cham: Springer International Publishing.
51. Prasad, B. V. V. S., Mandapati, S., Haritha, B., & Begum, M. J. (2020, August). Enhanced Security for the authentication of Digital Signature from the key generated by the CSTRNG method. In *2020 Third International Conference on Smart Systems and Inventive Technology (ICSSIT)* (pp. 1088-1093). IEEE.
52. Mukiri, R. R., Kumar, B. S., & Prasad, B. V. V. (2019, February). Effective Data Collaborative Strain Using RecTree Algorithm. In *Proceedings of International Conference on Sustainable Computing in Science, Technology and Management (SUSCOM)*, Amity University Rajasthan, Jaipur-India.
53. Balaraju, J., Raj, M. G., & Murthy, C. S. (2019). Fuzzy-FMEA risk evaluation approach for LHD machine—A case study. *Journal of Sustainable Mining*, 18(4), 257-268.
54. Thirumoorthi, P., Deepika, S., & Yadaiah, N. (2014, March). Solar energy based dynamic sag compensator. In *2014 International Conference on Green Computing Communication and Electrical Engineering (ICGCCEE)* (pp. 1-6). IEEE.
55. Vinayasree, P., & Reddy, A. M. (2025). A Reliable and Secure Permissioned Blockchain-Assisted Data Transfer Mechanism in Healthcare-Based Cyber-Physical Systems. *Concurrency and Computation: Practice and Experience*, 37(3), e8378.
56. Acharjee, P. B., Kumar, M., Krishna, G., Raminenei, K., Ibrahim, R. K., & Alazzam, M. B. (2023, May). Securing International Law Against Cyber Attacks through Blockchain Integration. In *2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE)* (pp. 2676-2681). IEEE.
57. Ramineni, K., Reddy, L. K. K., Ramana, T. V., & Rajesh, V. (2023, July). Classification of Skin Cancer Using Integrated Methodology. In *International Conference on Data Science and Applications* (pp. 105-118). Singapore: Springer Nature Singapore.
58. LAASSIRI, J., EL HAJJI, S. A. I. D., BOUHDADI, M., AOUDE, M. A., JAGADISH, H. P., LOHIT, M. K., ... & KHOLLADI, M. (2010). Specifying Behavioral Concepts by engineering language of RM-ODP. *Journal of Theoretical and Applied Information Technology*, 15(1).
59. Prasad, D. V. R., & Mohanji, Y. K. V. (2021). FACE RECOGNITION-BASED LECTURE ATTENDANCE SYSTEM: A SURVEY PAPER. *Elementary Education Online*, 20(4), 1245-1245.
60. Dasu, V. R. P., & Gujjari, B. (2015). Technology-Enhanced Learning Through ICT Tools Using Aakash Tablet. In *Proceedings of the International Conference on Transformations in Engineering Education: ICTIEE 2014* (pp. 203-216). Springer India.
61. Reddy, A. M., Reddy, K. S., Jayaram, M., Venkata Maha Lakshmi, N., Aluvalu, R., Mahesh, T. R., ... & Stalin Alex, D. (2022). An efficient multilevel thresholding scheme for heart image segmentation using a hybrid generalized adversarial network. *Journal of Sensors*, 2022(1), 4093658.
62. Srinivasa Reddy, K., Suneela, B., Inthiyaz, S., Hasane Ahammad, S., Kumar, G. N. S., & Mallikarjuna Reddy, A. (2019). Texture filtration module under stabilization via random forest optimization methodology. *International Journal of Advanced Trends in Computer Science and Engineering*, 8(3), 458-469.
63. Ramakrishna, C., Kumar, G. K., Reddy, A. M., & Ravi, P. (2018). A Survey on various IoT Attacks and its Countermeasures. *International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)*, 5(4), 143-150.
64. Sirisha, G., & Reddy, A. M. (2018, September). Smart healthcare analysis and therapy for voice disorder using cloud and edge computing. In *2018 4th international conference on applied and theoretical computing and communication technology (iCATccT)* (pp. 103-106). IEEE.
65. Reddy, A. M., Yarlagadda, S., & Akkinen, H. (2021). An extensive analytical approach on human resources using random forest algorithm. *arXiv preprint arXiv:2105.07855*.
66. Kumar, G. N., Bhavanam, S. N., & Midasala, V. (2014). Image Hiding in a Video-based on DWT & LSB Algorithm. In *ICPVS Conference*.
67. Naveen Kumar, G. S., & Reddy, V. S. K. (2022). High performance algorithm for content-based video retrieval using multiple features. In *Intelligent Systems and Sustainable Computing: Proceedings of ICISCC 2021* (pp. 637-646). Singapore: Springer Nature Singapore.
68. Reddy, P. S., Kumar, G. N., Ritish, B., SaiSwetha, C., & Abhilash, K. B. (2013). Intelligent parking space

- detection system based on image segmentation. *Int J Sci Res Dev*, 1(6), 1310-1312.
69. Naveen Kumar, G. S., Reddy, V. S. K., & Kumar, S. S. (2018). High-performance video retrieval based on spatio-temporal features. *Microelectronics, Electromagnetics and Telecommunications*, 433-441.
 70. Kumar, G. N., & Reddy, M. A. BWT & LSB algorithm based hiding an image into a video. *IJESAT*, 170-174.
 71. Lopez, S., Sarada, V., Praveen, R. V. S., Pandey, A., Khuntia, M., & Haralayya, D. B. (2024). Artificial intelligence challenges and role for sustainable education in india: Problems and prospects. *Sandeep Lopez, Vani Sarada, RVS Praveen, Anita Pandey, Monalisa Khuntia, Bhadrappa Haralayya (2024) Artificial Intelligence Challenges and Role for Sustainable Education in India: Problems and Prospects. Library Progress International*, 44(3), 18261-18271.
 72. Yamuna, V., Praveen, R. V. S., Sathya, R., Dhivva, M., Lidiya, R., & Sowmiya, P. (2024, October). Integrating AI for Improved Brain Tumor Detection and Classification. In *2024 4th International Conference on Sustainable Expert Systems (ICSES)* (pp. 1603-1609). IEEE.
 73. Kumar, N., Kurkute, S. L., Kalpana, V., Karuppanan, A., Praveen, R. V. S., & Mishra, S. (2024, August). Modelling and Evaluation of Li-ion Battery Performance Based on the Electric Vehicle Tiled Tests using Kalman Filter-GBDT Approach. In *2024 International Conference on Intelligent Algorithms for Computational Intelligence Systems (IACIS)* (pp. 1-6). IEEE.
 74. Sharma, S., Vij, S., Praveen, R. V. S., Srinivasan, S., Yadav, D. K., & VS, R. K. (2024, October). Stress Prediction in Higher Education Students Using Psychometric Assessments and AOA-CNN-XGBoost Models. In *2024 4th International Conference on Sustainable Expert Systems (ICSES)* (pp. 1631-1636). IEEE.
 75. Anuprathibha, T., Praveen, R. V. S., Sukumar, P., Suganthi, G., & Ravichandran, T. (2024, October). Enhancing Fake Review Detection: A Hierarchical Graph Attention Network Approach Using Text and Ratings. In *2024 Global Conference on Communications and Information Technologies (GCCIT)* (pp. 1-5). IEEE.
 76. Shinkar, A. R., Joshi, D., Praveen, R. V. S., Rajesh, Y., & Singh, D. (2024, December). Intelligent solar energy harvesting and management in IoT nodes using deep self-organizing maps. In *2024 International Conference on Emerging Research in Computational Science (ICERCS)* (pp. 1-6). IEEE.
 77. Praveen, R. V. S., Hemavathi, U., Sathya, R., Siddiq, A. A., Sanjay, M. G., & Gowdish, S. (2024, October). AI Powered Plant Identification and Plant Disease Classification System. In *2024 4th International Conference on Sustainable Expert Systems (ICSES)* (pp. 1610-1616). IEEE.
 78. Dhivya, R., Sagili, S. R., Praveen, R. V. S., VamsiLala, P. N. V., Sangeetha, A., & Suchithra, B. (2024, December). Predictive Modelling of Osteoporosis using Machine Learning Algorithms. In *2024 4th International Conference on Ubiquitous Computing and Intelligent Information Systems (ICUIS)* (pp. 997-1002). IEEE.
 79. Kemmannu, P. K., Praveen, R. V. S., Saravanan, B., Amshavalli, M., & Banupriya, V. (2024, December). Enhancing Sustainable Agriculture Through Smart Architecture: An Adaptive Neuro-Fuzzy Inference System with XGBoost Model. In *2024 International Conference on Sustainable Communication Networks and Application (ICSCNA)* (pp. 724-730). IEEE.
 80. Praveen, R. V. S. (2024). *Data Engineering for Modern Applications*. Addition Publishing House.