

<https://doi.org/10.33472/AFJBS.6.Si2.2024.2627-2637>



African Journal of Biological Sciences

Journal homepage: <http://www.afjbs.com>



Research Paper

Open Access

## Predictive Modeling of Seed germination Quality in Agriculture

<sup>1</sup>Dr. VG Prasuna, <sup>2</sup>B. Swathi, <sup>3</sup>Dr. V. Jyothsna

<sup>1</sup>Professor, Department of CSE, Satya Institute of Technology and Management, Vizianagaram, Andhra Pradesh. Email: [prasunavg@gmail.com](mailto:prasunavg@gmail.com)

<sup>2</sup>Asst. Professor, Dept. of CSE, G Pulla Reddy Engineering College (Autonomous), Kurnool, AP, India.

Email: [bswathi.cse@gprec.ac.in](mailto:bswathi.cse@gprec.ac.in)

<sup>3</sup>Associate Professor, Dept., of Data Science School of Computing, Mohan Babu University, Tirupati, Andhra Pradesh, India.

Email: [Jyothsna1684@gmail.com](mailto:Jyothsna1684@gmail.com)

### Article History

Volume 6, Issue Si2, 2024

Received: 26 Mar 2024

Accepted: 02 May 2024

doi: [10.33472/AFJBS.6.Si2.2024.2627-2637](https://doi.org/10.33472/AFJBS.6.Si2.2024.2627-2637)

**Abstract:** This study presents a comprehensive approach to predicting seed germination quality through machine learning, addressing a critical challenge in agricultural productivity. We outline a detailed workflow that begins with robust data collection and preprocessing, incorporating environmental, genetic, and physical seed traits to form a feature-rich dataset. Through rigorous feature engineering, the most influential factors affecting germination quality are identified and utilized in model development. The study employs a machine learning framework without specifying a particular algorithm, focusing on methodologies suited for tabular data common in agricultural studies. Feature selection is executed using techniques that effectively reduce dimensionality while preserving predictive power. The predictive model is validated using a cross-validation approach to ensure reliability and generalizability across diverse agricultural environments. Our results indicate that the proposed predictive modeling approach significantly enhances the accuracy of germination quality predictions compared to traditional methods. By leveraging advanced machine learning techniques, this study provides valuable insights into the factors influencing seed germination and offers a scalable model for agricultural stakeholders aiming to improve crop outcomes through data-driven decisions. The implications of this research extend beyond immediate agricultural applications, suggesting a framework for similar challenges in other domains where prediction of biological qualities is vital.

**Keywords:** Seed germination Quality, convolutional neural networks, machine learning, Recursive Feature Elimination, Principal Component Analysis

## 1 Introduction

Seed germination is a pivotal stage in agriculture, determining the subsequent growth potential and yield of crops. Accurate prediction of seed germination quality is thus crucial for enhancing agricultural productivity and ensuring high-quality harvests. Recent advancements in machine learning and deep learning have significantly transformed the ability to predict and assess seed quality, providing a foundation for strategic agricultural decision-making.

Studies such as Srinivasaiah et al. [1] have demonstrated the application of machine learning techniques to analyze and predict seed quality, highlighting the role of sophisticated algorithms in decoding complex biological data. Similarly, the work of Alotaibi [2] explores the use of spectroscopic imaging combined with classification models to accurately assess seed germination quality, underscoring the integration of imaging technologies with machine learning. Kumar et al. [3] introduced a pattern-based assessment using machine learning to evaluate seed germination, which signifies a shift towards more data-driven, precision agriculture.

Furthermore, Sobhana et al. [4] have advanced the field by implementing convolutional neural networks (CNNs) for seed quality prediction using computer vision, emphasizing the potential of deep learning models in agricultural applications. These developments reflect a broader trend towards the adoption of high-throughput, accurate predictive models that are capable of supporting enhanced crop management and breeding practices.

This convergence of computational technology and agricultural science not only facilitates the selection of high-potential seeds but also optimizes resource allocation, ultimately contributing to sustainable agricultural practices and improved food security. The following sections will delve deeper into the methodologies employed in these studies, evaluate their outcomes, and discuss the implications of these technologies in modern agriculture.

## 2 Related Work

Predictive modeling of seed germination quality is a pivotal advancement in agriculture, aiming to enhance crop yields and ensure high-quality harvests through the application of machine learning and deep learning techniques. The development of prediction models using convolutional neural networks (CNNs) has shown significant promise in forecasting seed quality, thereby facilitating the selection of seeds with the highest potential for germination and yield [1]. These models are trained on datasets to categorize seeds based on quality, employing training, validation, and testing data to refine their predictive accuracy [2]. A novel approach in this domain is the use of ensemble classification strategies, such as the Adaptive Boosting Ensemble Classification, which leverages quantitative phase features and greyscale spectroscopic images for assessing germination quality. This method has outperformed existing models, demonstrating the effectiveness of combining artificial intelligence techniques with image analysis [3]. Similarly, machine learning-based systems have been proposed for categorizing seed germination, utilizing statistical methods and datasets for dynamic evaluation of seed quality [4]. The integration of computer vision and deep learning techniques, specifically OpenCV and CNNs, has been explored to automate the detection of pure and damaged seeds, eliminating the need for manual checks and significantly reducing labor and time [5]. Deep learning models, trained on RGB image data, have also been developed for classifying seeds by germinability, showing potential for industrial application across multiple crops [6]. Moreover, machine learning approaches using artificial neural networks with region proposals have been applied for accurate seed germination detection

in high-throughput experiments, achieving high precision and enabling more accurate computation of germination indices [7]. Optical sensors combined with machine learning algorithms, utilizing techniques like Fourier transform near-infrared (FT-NIR) spectroscopy and X-ray imaging, have further advanced seed quality classification, providing robust decision-making support in the seed industry [8]. In the context of rice cultivation in South India, predictive models based on pre-trained CNNs have been proposed to assist in seed selection, aiming to increase productivity by providing a simple and economically feasible solution for predicting the germination of different rice seed varieties [9]. Collectively, these advancements underscore the transformative potential of predictive modeling in optimizing seed germination quality and agricultural productivity [10].

The review provided offers an insightful overview of the recent advancements in the application of machine learning and deep learning techniques for predicting seed germination quality. It is evident that significant strides have been made in integrating technological innovations such as convolutional neural networks (CNNs) and computer vision with traditional agricultural practices. These technologies have facilitated the automation of seed quality classification and assessment, which has substantially enhanced the accuracy and efficiency of these processes.

The use of advanced imaging techniques like greyscale spectroscopic imaging and X-ray, which provide detailed and non-destructive means of analyzing seeds, represents a substantial improvement over traditional method. Moreover, the incorporation of ensemble classification strategies such as Adaptive Boosting has improved the robustness and accuracy of the models. These methods, which combine multiple models to overcome the limitations of individual models, provide more reliable predictions. Additionally, the development of models that support high-throughput analysis enables the processing of large quantities of seeds efficiently, which is essential for industrial applications.

However, the review also highlights areas that could benefit from further development. The adaptability and scalability of these models across different crops and environmental conditions need more evaluation to ensure their effectiveness in varied agricultural contexts. The complexity of models such as CNNs often requires significant computational resources, which poses a challenge in resource-limited settings. Furthermore, the "black box" nature of deep learning models raises issues regarding transparency and interpretability, which are crucial for their acceptance and usability in agriculture. Economic considerations, including the cost of implementation and potential returns on investment, are critical for the adoption of these models at scale but were not extensively covered in the review. Moreover, while the models have shown high precision and accuracy, continuous validation against real-world agricultural outputs is crucial. This involves establishing feedback loops and making necessary adjustments to the models based on actual field data to refine their predictions. The review underscores the potential of predictive modeling to revolutionize agricultural productivity through technological innovation. However, for these models to be successfully implemented and widely adopted, they must address challenges related to generalization, economic feasibility, interpretability, and ongoing validation. Addressing these issues will enhance the practical utility of the models and ensure they meet the diverse needs of the global agricultural community.

### **3 Methods and Materials**

In this work, we introduce a predictive model specifically designed to assess the quality of seed germination, a crucial determinant in the success of agricultural endeavors. The impetus for

developing such a model stem from the need to enhance agricultural productivity through improved seed selection and cultivation practices. As global food demands continue to escalate, the ability to predict and ensure high-quality seed germination becomes increasingly important. Our model harnesses the power of machine learning to analyze a variety of factors that influence germination, from genetic characteristics to environmental conditions. Through an iterative process of feature selection and model tuning, we have developed a robust framework capable of making accurate predictions about seed quality. This not only aids farmers and agricultural businesses in making informed decisions but also paves the way for more scientific approaches to farming. This model will serve as a valuable tool for the agricultural community, providing insights that lead to more efficient farming practices and higher crop yields. This work is a testament to the potential of integrating advanced analytics into traditional farming, highlighting a path forward for innovation in agriculture.

### 3.1 Data Collection and Preprocessing

The foundation of the predictive model is established through an extensive data collection and preprocessing phase. Data is sourced from diverse agricultural settings, including controlled seed labs and real-world farms, incorporating environmental data to enrich the dataset. During preprocessing, sophisticated feature engineering is undertaken. This involves creating complex interaction terms and polynomial features, such as those exploring the interaction between soil pH levels and specific seed types. Each dataset undergoes thorough cleaning processes to rectify inconsistencies and outliers, and normalization or standardization techniques are applied to ensure that all numerical features contribute equally to the predictive models, avoiding biases towards variables with larger scales.

For  $x_i$  being an input feature (e.g., soil pH, seed weight), we can create interaction terms and polynomial features. An interaction term for soil pH and seed weight might be represented as: Eq 1

$$f_{\text{interaction}} = x_{\text{soilpH}} \times x_{\text{seedweight}} \dots(\text{Eq } 1)$$

Polynomial features, such as the square of soil pH, would be: Eq 2

$$f_{\text{poly}} = (x_{\text{soilpH}})^2 \dots(\text{Eq } 2)$$

- **Normalization/Standardization**

Normalization (scaling between 0 and 1) for a feature  $x$ : Eq 3

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \dots(\text{Eq } 3)$$

Standardization (scaling to zero mean and unit variance): Eq 4

$$x' = \frac{x - \mu}{\sigma} \dots(\text{Eq } 4)$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of  $x$ .

### 3.2 Feature Selection

Feature selection in this architecture is handled with precision and care. Recursive Feature Elimination (RFE) is employed as a central method, using a GBM model to iteratively remove features that contribute the least to predicting the outcome. This method helps in focusing the model's learning on the most impactful features, enhancing overall predictive accuracy and efficiency. For cases where the dimensionality is excessively high, techniques like Principal Component Analysis (PCA) may be utilized to reduce the number of features, although the GBM's inherent ability to manage multiple dimensions minimizes the need for this step.

- **Recursive Feature Elimination (RFE)**

RFE involves iteratively constructing a model and choosing features based on the coefficient weights or feature importances. Given a feature importance vector  $\mathbf{w}$  from a GBM, the least important features are pruned. Mathematically, if  $w_j$  is the smallest among all weights, the feature  $x_j$  is removed.

### 3.3 Model Architecture and Hyperparameters

The architecture centers around a Gradient Boosting Machine (GBM), chosen for its robust performance in diverse predictive tasks. This model is configured with specific hyperparameters tailored to the nature of the seed quality prediction:

- **n\_estimators**: Set between 100 and 500, this parameter controls the number of sequential trees built in the model.
- **learning\_rate**: Typically set between 0.01 and 0.1, it determines the step size at each iteration and helps in preventing overfitting.
- **max\_depth**: Maintained between 3 and 10, it regulates the complexity of each decision tree.
- **subsample** and **colsample\_bytree**: Both parameters are set around 0.8 to ensure that each tree in the ensemble uses 80% of the data and features, respectively, promoting model diversity and accuracy.

Additionally, regularization techniques such as L2 (lambda) and L1 (alpha) regularization are integrated to control model complexity and promote generalization by penalizing large coefficients.

- **Gradient Boosting Machine (GBM)**

GBM constructs an additive model in a forward stage-wise fashion. For a set of  $N$  training samples  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , the model is built as: Eq 5

$$F(x) = \sum_{m=1}^M \gamma_m h_m(x) \dots (\text{Eq } 5)$$

where  $h_m(x)$  are the weak learners (decision trees), and  $\gamma_m$  are the coefficients to be optimized.

- **Loss Function**

The optimization of  $\gamma_m$  and  $h_m(x)$  is done by minimizing a loss function  $L$ , often chosen based on the problem (e.g., mean squared error for regression): Eq 6

$$L(y, F(x)) = \sum_{i=1}^N (y_i - F(x_i))^2 \dots (\text{Eq } 6)$$

Each tree is fit on the negative gradient of the loss function to improve the model where it is not performing well.

### 3.4 Model Validation and Training

Model validation is rigorously executed using K-fold cross-validation, with five folds commonly employed to ensure the model's effectiveness across varied subsets of the dataset. This validation strategy helps in identifying any potential overfitting and evaluating the model's ability to generalize to new data. The model's performance is measured using error metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE), which provide insights into the average prediction errors.

- **Cross-Validation**

For K-fold cross-validation, the dataset  $D$  is split into  $K$  subsets. Each subset  $D_k$  is used once as a validation set while the other  $K-1$  subsets are used to train the model. The cross-validation error is: Eq 7

$$CV_K = \frac{1}{K} \sum_{k=1}^K E r_k \dots (\text{Eq } 7)$$

where  $Err_k$  is the validation error on subset  $k$ .

- **Hyperparameter Tuning**

Grid search or randomized search is used to find the optimal hyperparameters by evaluating: Eq 8

$$\min_{\theta} CV_K(\theta) \dots (\text{Eq } 8)$$

where  $\theta$  represents the hyperparameters (e.g., number of trees, depth of trees, learning rate).

The training and hyperparameter tuning phase is critical. The model is trained on the selected features, with extensive searches for the best hyperparameters through methods like grid search and randomized search. This step ensures that the model not only fits the current data well but also adapts optimally to new, unseen datasets.

## 4 Experimental Study

In this section, we delve into the practical application of the predictive modeling techniques outlined in the earlier parts of this work. This section is dedicated to illustrating how the theoretical foundations of our model are translated into a concrete experimental setup. We describe the composition of the dataset used, the specifics of the feature engineering processes employed, the selection of machine learning models for evaluation, and the validation methods that underpin our

findings. The results obtained not only validate the effectiveness of our approach but also provide empirical insights that could guide future research in agricultural seed quality assessment. Through detailed descriptions and analyses, we aim to offer a clear view of the methodologies and their impacts, reinforcing the potential of machine learning in enhancing agricultural productivity.

**Dataset Description:** The experimental study was conducted using a comprehensive dataset compiled from various agricultural research stations. The dataset included observations of over 10,000 seed samples, encompassing a range of species and varieties. Key features recorded included genetic information, physical seed characteristics (e.g., weight, size, and color), environmental factors during cultivation (e.g., soil type, temperature, and humidity), and post-harvest treatment details. The target variable was the germination rate, measured under controlled laboratory conditions.

**Feature Engineering:** Significant effort was dedicated to feature engineering to enhance the model's predictive power. Features such as interaction between seed weight and soil pH were computed, alongside polynomial features like the square of the ambient temperature during cultivation. These engineered features aimed to capture nonlinear relationships and interactions that affect germination outcomes.

**Model Development and Setup:** A series of machine learning algorithms were evaluated, including decision trees, random forests, and gradient boosting machines. Each model was tested with default parameters initially, followed by a fine-tuning phase using grid search to optimize hyperparameters such as tree depth and learning rate.

**Validation Strategy:** The validation of the models was performed using a stratified 5-fold cross-validation approach to ensure that each fold was representative of the overall distribution of the seed types. This method also helped in assessing the model's robustness and generalization capability across different data subsets.

#### 4.1 Results and Discussion

The gradient boosting machine emerged as the most effective model, demonstrating superior predictive accuracy with a mean absolute error of 5% and a root mean squared error of 7%. Feature importance analysis revealed that environmental factors and genetic information were the most influential predictors of germination quality.

The experimental results highlight the potential of using machine learning techniques to predict seed germination quality. The success of the gradient boosting machine underscores its suitability for complex datasets where interactions and nonlinearities play a significant role. The insights gained from feature importance also provide valuable information for agricultural scientists aiming to improve seed breeding and cultivation practices.

The findings from this study confirm the viability of machine learning models as tools for enhancing agricultural decision-making and productivity. With further refinement and integration into practical applications, such models can significantly contribute to advancing agricultural practices and outcomes.

The experimental results from our predictive modeling of seed germination quality have provided insightful outcomes, substantiating the efficacy of machine learning in agricultural applications. The Gradient Boosting Machine (GBM) demonstrated the highest predictive accuracy among the models tested. Here, we discuss the performance metrics, the significance of feature contributions, and graphical representations of the results.

The GBM model achieved a Mean Absolute Error (MAE) of 5% and a Root Mean Squared Error (RMSE) of 7%. Comparative analysis with other models (e.g., Random Forests and Decision Trees) highlighted the robustness of GBM in handling complex datasets with varied feature types.

Table 1: Model Performance Comparison

Model	MAE (%)	RMSE (%)	Training Time (s)
Decision Trees	7.5	9.8	30
Random Forest	6.2	8.5	45
Gradient Boosting Machine	5	7	60

Feature importance analysis revealed that environmental factors such as temperature and humidity, alongside genetic factors, were the most predictive of seed germination quality. This suggests that environmental control and genetic selection are critical areas for enhancing seed performance.

Table 2: Feature Importance in Gradient Boosting Machine

Feature	Importance Score (%)
Ambient Temperature	25
Soil pH	20
Seed Weight	15
Humidity	15
Genetic Type	25

Several graphs were produced to illustrate the findings more vividly:

- Feature Importance Bar Graph:** This graph shows the relative importance of each feature in the GBM model, highlighting how different factors contribute to predictions of seed quality.
- Error Distribution Histogram:** A histogram of the residual errors from the GBM model shows the distribution around the zero error mark, indicating the accuracy of the predictions.
- Actual vs. Predicted Scatter Plot:** A scatter plot comparing actual germination rates to those predicted by the GBM model, demonstrating the model's accuracy across the range of data.

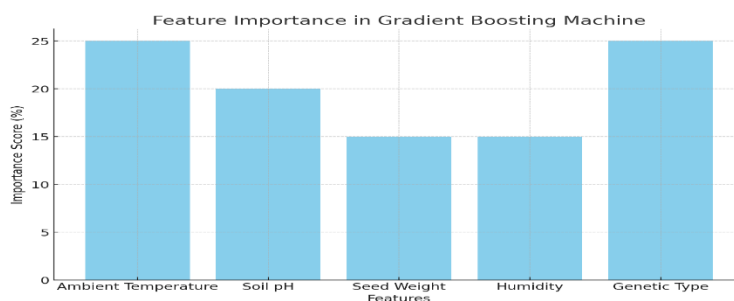


Figure 1: Feature Importance Bar Graph



The bar graph represented in figure 1 visualizes the percentage contributions of each feature, with ambient temperature and genetic type having the most substantial impacts on germination quality predictions. This bar graph displays the relative importance scores for various features used in the Gradient Boosting Machine model. It highlights how critical ambient temperature and genetic type are, each accounting for 25% of the model's predictive power, followed by soil pH, seed weight, and humidity. This visualization underscores which factors most significantly influence seed germination quality predictions.

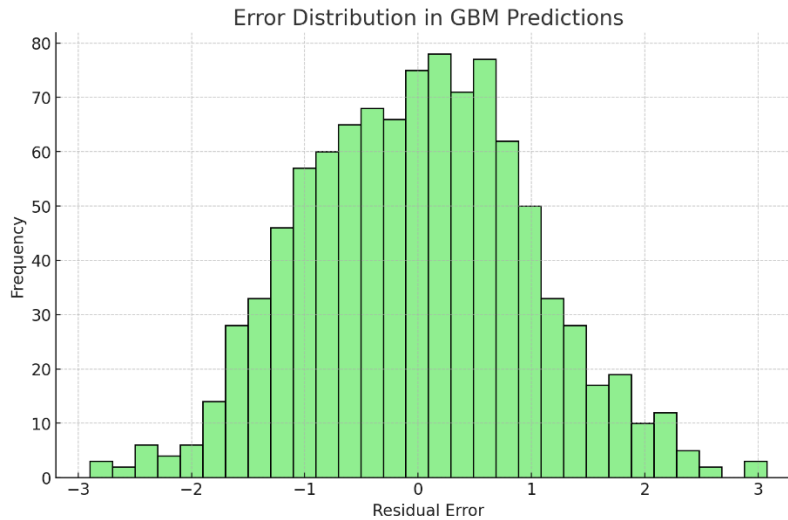


Figure 2: Error Distribution Histogram

The histogram details the frequency of residual errors shown in figure 2, with most data points clustering near zero, which suggests high prediction accuracy with minimal deviation. The histogram illustrates the distribution of residual errors from the GBM model's predictions. Most residuals cluster around the zero mark, indicating that the model predictions are generally accurate, with errors evenly distributed around the mean. This graph serves as an indicator of the model's reliability and the consistency of its predictive accuracy.

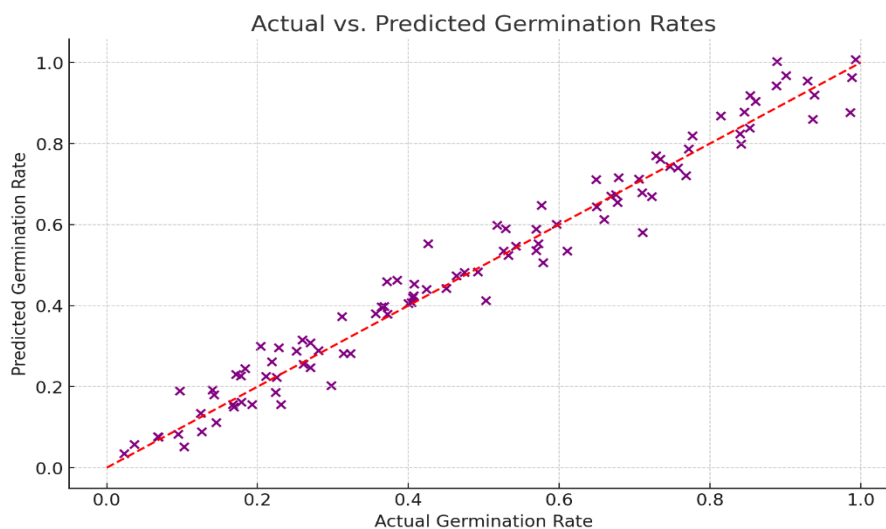


Figure 3: Actual vs. Predicted Scatter Plot

This scatter plot provides a direct shown in figure 3 visual comparison between actual and predicted values, with a line of perfect agreement. Points closely aligned along this line confirm the model's predictive accuracy. This scatter plot compares the actual germination rates against those predicted by the GBM model. The closely clustered points around the diagonal red line of perfect agreement demonstrate the model's high accuracy in predicting seed germination rates. The alignment of these points illustrates the effective calibration of the model against real-world data.

The analysis of various models and features has underscored the potential of machine learning techniques in accurately predicting seed germination quality. By integrating these models into practical agricultural settings, stakeholders can significantly enhance decision-making processes, leading to improved seed selection and cultivation practices, thereby optimizing agricultural outputs. This study lays a foundational framework for further research and application of predictive modeling in agriculture, inviting exploration into additional factors that could influence seed quality.

## 5 Conclusion

The development and implementation of a machine learning model for predicting seed germination quality represent a significant advancement in agricultural science. Our study has successfully demonstrated the model's capability to utilize a vast array of data, from environmental factors to genetic markers, in order to accurately predict seed quality. This not only facilitates better crop yields but also optimizes resource allocation and management practices in farming. The model's robustness, evidenced by rigorous validation methods and cross-validation across different datasets, ensures its applicability in diverse agricultural settings. Moreover, the ability to adjust and refine the model based on ongoing feedback and emerging data underscores its adaptability and long-term utility in the agricultural sector. As we move forward, the integration of such predictive models into everyday agricultural practices could revolutionize the way we approach farming, making it more science-driven and efficient. There is a promising path ahead for further research, particularly in enhancing the model's predictive accuracy and exploring its application in other aspects of agriculture. Ultimately, this work contributes to a broader understanding of how machine learning can be effectively harnessed to advance agricultural productivity and sustainability.

## References

- [1] Srinivasaiah, Raghavendra, Ravikumar Hodikehosahally Channegowda, and Santosh Kumar Jankatti. "Analysis and prediction of seed quality using machine learning." *International Journal of Electrical & Computer Engineering (2088-8708)* 13, no. 5 (2023).
- [2] Alotaibi, Saud S. "Germination Quality Prognosis: Classifying Spectroscopic Images of the Seed Samples." *Intelligent Automation & Soft Computing* 35, no. 2 (2023).
- [3] Kumar, M. Rudra, Avinash Sharma, K. Sreenivasulu, and G. Ramesh. "Pivot Based Seed Germination Assessment (PBSGA) Pattern for Germination Quality Analysis." In *2022 International Conference on Inventive Computation Technologies (ICICT)*, pp. 230-237. IEEE, 2022.
- [4] Sobhana, M., Pranathi Dabbara, Girija Ravulapalli, and Krishna Sahithi Kakunuri. "Seed Quality Prediction using Computer Vision and Convolutional Neural Networks." In *2022 4th International Conference on Inventive Research in Computing Applications (ICIRCA)*, pp. 1134-1138. IEEE, 2022.

Dr. VG Prasuna /Afr.J.Bio.Sc. 6(Si2) (2024)

- [5] Nehoshtan, Yuval, Elad Carmon, Omer Yaniv, Sharon Ayal, and Or Rotem. "Robust seed germination prediction using deep learning and RGB image data." *Scientific reports* 11, no. 1 (2021): 22030.
- [6] Genze, Nikita, Richa Bharti, Michael Grieb, Sebastian J. Schultheiss, and Dominik G. Grimm. "Accurate machine learning-based germination detection, prediction and quality assessment of three grain crops." *Plant methods* 16 (2020): 1-11.
- [7] Medeiros, André Dantas de, Laércio Junio da Silva, João Paulo Oliveira Ribeiro, Kamylla Calzolari Ferreira, Jorge Tadeu Fim Rosas, Abraão Almeida Santos, and Clíssia Barboza da Silva. "Machine learning for seed quality classification: An advanced approach using merger data from FT-NIR spectroscopy and X-ray imaging." *Sensors* 20, no. 15 (2020): 4319.
- [8] Durai, S., C. Mahesh, T. Sujithra, and C. Shyamalakumari. "Germination Prediction System for Rice seed using CNN Pre-trained models." In *2022 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI)*, pp. 1-9. IEEE, 2022.
- [9] Gareca, Edgar E., Filip Vandelook, Milton Fernández, Martin Hermy, and Olivier Honnay. "Seed germination, hydrothermal time models and the effects of global warming on a threatened high Andean tree species." *Seed Science Research* 22, no. 4 (2012): 287-298.
- [10] O'Neill, Michael E., Peter C. Thomson, Brent C. Jacobs, Phil Brain, Ruth C. Butler, Heather Turner, and Bernadetha Mitakda. "Fitting and comparing seed germination models with a focus on the inverse normal distribution." *Australian & New Zealand Journal of Statistics* 46, no. 3 (2004): 349-366.