

AN INTUITIVE ANALYSIS ON EARLY DETECTION WITH IAGT MODEL FOR COTTON CROP YIELD PREDICTION

PORANDLA SRINIVAS¹, DR. SURESH A²

¹Research Scholar, Department of Networking and Communications, School of Computing, SRM Institute of Science and Technology, Kattankulathur, Chennai, Tamil Nadu, India

²Associate Professor, Department of Networking and Communications, School of Computing, SRM Institute of Science and Technology, Kattankulathur, Chennai, Tamil Nadu, India

E-mail: ¹srinivas.research@yahoo.com, ²prisu6esh@yahoo.com

ABSTRACT

Cotton yield prediction plays a critical role in modern agriculture, influencing food security, economic stability, and effective resource management. Despite the advancements in various predictive algorithms like support vector machines (SVM), random forests, and artificial neural networks, challenges persist in handling the complex, high-dimensional, and non-linear nature of agricultural data. Traditional models struggle with issues such as dynamic environmental fluctuations, incomplete datasets, and an inability to effectively adapt to evolving conditions in the field. As a result, these models often fail to provide accurate and timely predictions, leading to inefficiencies in crop management, resource allocation, and risk mitigation. This study addresses the knowledge gap in cotton yield prediction by introducing the Integrated Adaptive Growth Tree (IAGT) algorithm, which combines decision trees with deep learning techniques for real-time adaptation to changing agricultural conditions. By integrating multi-source data, including satellite imagery, weather forecasts, and soil sensor readings, the IAGT model offers a novel approach to yield forecasting, surpassing traditional methods in both accuracy and adaptability. Simulations using both original and synthetic cotton yield datasets showed a remarkable 98% accuracy, demonstrating significant improvements in prediction performance. This study not only provides new insights into the effective integration of diverse data sources for crop yield forecasting but also introduces a robust framework for early-stage disease detection and anomaly identification, thus contributing to the growing field of precision agriculture. The IAGT model's success in enhancing cotton yield predictions sets the stage for broader applications in crop management and agricultural sustainability.

Keywords: *Integrated Adaptive Growth Tree (IAGT), Convolution Neural Networks (CNN), Gated Adversarial Neural Networks (GAN), Support Vector machines (SVM),*

1. INTRODUCTION

Cotton is a vital crop globally, with India and the United States being two of the largest producers. In India, cotton farming plays a crucial role in the economy, providing employment to millions of farmers and supporting the textile industry. Cotton production is especially important in states like Maharashtra, Gujarat, and Telangana. Despite India being the world's largest producer of cotton, predicting cotton yields accurately remains a challenge due to the diverse climatic conditions, soil types, and variations in farming practices across the country. Accurate cotton yield predictions are essential for optimizing agricultural practices, improving crop management, reducing resource wastage, and ensuring food security. Similarly, in the United States, cotton farming is concentrated in states like Texas, Georgia, and Mississippi, where it

is a major contributor to the agricultural economy. While U.S. cotton farmers have increasingly adopted precision agriculture technologies, yield prediction continues to be an area requiring improvement, especially with the challenges posed by climate change and pest outbreaks. Models that can predict cotton yield with high accuracy are needed to assist farmers in decision-making and to improve the overall efficiency of cotton farming in both countries [1].

1.1 Existing Algorithms and Their Gaps

Various algorithms, such as machine learning models, regression analysis, and deep learning approaches, have been employed for cotton yield prediction. For instance, methods like Support Vector Machines (SVM), Random Forest, and

Artificial Neural Networks (ANNs) have been explored for their ability to model complex non-linear relationships in agricultural data [6],[7]. Other approaches use satellite-based data, weather forecasts, and remote sensing technologies to predict crop yields [2]. However, these models often face limitations in accurately capturing the full spectrum of factors that influence cotton yield, such as micro-climatic variations, soil health, pest pressure, and irrigation levels. Additionally, existing models struggle with scalability, especially in regions with highly variable environmental conditions. There is also a lack of real-time data integration, which makes these models less adaptive to changing weather patterns and farming practices. Despite these advancements, cotton yield prediction continues to suffer from inaccuracies, especially in regions with unpredictable weather and pest dynamics.

1.2 Problem Statement

Previous studies on cotton yield prediction have explored a range of approaches, each contributing to the overall understanding of the factors that influence yield outcomes. The study in [2] highlighted the integration of Earth observation variables with machine learning models, achieving high accuracy, though they noted limitations in real-time prediction adjustments and reliance on high-quality satellite data. Similarly, the usage of high-resolution spatial data to reveal cotton yield variations, but their approach lacked integration with other crucial environmental factors, such as pest infestations and diseases are implemented in [3]. The authors in [6] focused on an ensemble machine learning model using weather parameters, showing improved accuracy but encountering challenges due to inconsistent weather data quality across regions. Meanwhile, the work in [7] has proposed a machine learning-based approach combining field and synthetic data, yet their method faced scalability issues in diverse farming regions. Other works, like those by [9], explored hardware acceleration and disease classification through deep learning, respectively, yet they did not address the integration of disease data with yield prediction models. Finally, the work stated in [10] has examined the potential of multispectral and thermal sensors, contributing to more accurate cotton production estimates, though their study called for better adaptation to various environmental conditions.

While these studies have made valuable strides in cotton yield prediction, many still face challenges in terms of scalability, real-time adaptability, and integration with dynamic environmental factors. The

novelty of our work lies in its ability to address these gaps by introducing the Integrated Adaptive Growth Tree (IAGT) algorithm, which combines decision trees with deep learning techniques for real-time adaptation to fluctuating agricultural conditions. Our model integrates diverse data sources such as satellite imagery, weather forecasts, soil sensors, and disease information, offering a more holistic and adaptive approach. Unlike previous models that focus solely on one type of data or environmental factor, our approach provides a comprehensive solution, improving prediction accuracy to 98% while also enabling early-stage disease detection and anomaly identification. This makes our work distinct in its motivation to bridge the existing gaps and in its findings, which offer a more robust and practical framework for precision agriculture.

1.3 Introduction of the Proposed IAGT Algorithm

The Integrated Adaptive Growth Tree (IAGT) algorithm presents a novel solution to these challenges. By combining the strengths of decision trees with deep learning, IAGT offers a more flexible and adaptive approach to cotton yield prediction. The decision tree component helps handle non-linear relationships between various environmental and agricultural parameters, such as temperature, rainfall, and soil quality, while deep learning enhances the model's ability to learn from large datasets, enabling better predictions based on historical and real-time data. Moreover, the adaptive nature of IAGT allows it to continuously adjust to new data, making it more resilient to the dynamic nature of farming conditions and improving its predictive accuracy. The IAGT algorithm can integrate multi-source data, including weather data, soil health indicators, satellite imagery, and even farmer inputs, providing a more holistic view of factors affecting cotton yield. This model is designed to be robust, scalable, and adaptable to various agricultural zones, including India and the U.S., where cotton farming faces distinct challenges.

1.4 Simulations and Deep Learning Metrics for Cotton Yield Prediction

In simulations using both real-world and synthetic data, the proposed IAGT algorithm demonstrated its capability to predict cotton yields with up to 98% accuracy. By training the model with large datasets collected from multiple regions, it was able to learn complex patterns and relationships that traditional

models failed to capture. For instance, the inclusion of real-time weather data, soil health indicators, and pest management information enhanced the model's predictive power. The model's success can be attributed to its ability to adapt to different data distributions and environmental factors, which are crucial for cotton farming. Metrics such as Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and R^2 values were used to evaluate the model's performance. The results showed that the IAGT algorithm outperformed conventional models, making it a promising tool for future cotton yield predictions. With this high level of accuracy, the algorithm can be used to inform decision-making for farmers, optimize resource usage, and improve the sustainability of cotton farming, especially in regions affected by climate change and inconsistent weather patterns.

1.5 Contribution to Society and Future Implications

The integration of the IAGT algorithm into cotton farming will not only provide more accurate yield predictions but also enable farmers to make informed decisions on crop management. This could lead to better resource allocation, such as the efficient use of water, fertilizers, and pesticides, reducing costs and environmental impact. Moreover, by offering real-time predictions, the algorithm can help farmers mitigate risks associated with pest outbreaks, crop diseases, and adverse weather conditions, ensuring stable yields and increasing farmers' resilience to climate change. With a focus on societal benefits, this technology can enhance food security and boost economic stability in cotton-growing regions, particularly in developing countries like India. By offering a more adaptive, scalable, and accurate solution for cotton yield prediction, the IAGT algorithm will play a key role in advancing sustainable agricultural practices, improving productivity, and contributing to the global cotton supply chain.

1.5 Objectives;

- **Precision Agriculture and Sustainability:** The IAGT algorithm enhances cotton yield predictions, helping farmers optimize resources like water, fertilizers, and pesticides, thus reducing waste and supporting sustainable farming practices.
- **Improved Prediction Accuracy:** By integrating multiple data sources (weather, soil health, satellite imagery), the IAGT algorithm provides more accurate and

adaptive cotton yield forecasts, capturing complex relationships that traditional models overlook.

- **Economic Stability and Food Security:** The IAGT algorithm improves yield prediction accuracy, helping farmers plan better, reduce costs, and enhance productivity, leading to greater economic stability and food security in cotton-growing regions.

1.6 Overview of the Paper

The paper presents the Integrated Adaptive Growth Tree (IAGT) algorithm, a novel approach combining decision trees with deep learning for more accurate cotton yield prediction. The model effectively handles complex environmental data, including weather, soil health, and pest management, offering a more flexible and adaptive solution for cotton farming. The simulations show that the IAGT algorithm achieves up to 98% prediction accuracy, surpassing traditional models. The proposed model provides a scalable and robust framework that can be used across different agricultural zones, helping optimize resource usage, improve sustainability, and support decision-making in cotton farming. With its focus on integrating multi-source data, real-time adaptation, and enhanced predictive power, the IAGT algorithm represents a promising tool for advancing precision agriculture and ensuring more resilient farming practices.

2. LITERATURE SURVEY:

2.1 Introduction to Cotton Yield Prediction Models

The study of cotton yield prediction has evolved significantly with the integration of machine learning, optimization techniques, and advanced data analysis methods. Various approaches have been proposed to enhance the accuracy and reliability of cotton yield forecasts. In recent years, extended grey models based on particle swarm optimization have shown promising results in predicting cotton yield, providing valuable insights into crop management and production planning [1]. This approach combines the strengths of grey modeling with optimization algorithms, making it more effective in dealing with uncertainty in agricultural data. Another trend in cotton yield prediction focuses on the use of earth observation variables. By leveraging explainable boosting machines, researchers have explored how remote sensing data can improve the prediction of cotton

yield. These studies highlight the importance of incorporating diverse datasets, such as satellite imagery and environmental factors, to generate reliable predictions [2]. The ability to explain and interpret the machine learning models used for this purpose adds a layer of transparency that can assist farmers and researchers in making data-driven decisions.

2.2 Satellite Remote Sensing in Cotton Yield Prediction

The interpretability of cotton yield prediction models has become a focal point in recent research, with studies analyzing the influence of end-season and mid-season predictors on yield to optimize prediction accuracy. Time-sensitive variables, such as weather conditions and crop health, have proven critical in enhancing forecast precision, and the use of interpretable models builds trust in these predictions, which is essential for practical agricultural applications [3]. In line with this, the growing use of time-series data in agriculture has introduced a methodology focused on integrating remote sensing data over time. These interpretable models capture the dynamic nature of cotton growth and yield patterns, providing valuable insights for monitoring crop health and adjusting farming practices based on the predictions [4]. Furthermore, the incorporation of high-resolution data and statistical analysis techniques has gained attention, particularly through studies like the nationwide research in Turkey, which demonstrates how low-resolution statistics can uncover high-resolution yield variations, ultimately improving regional predictions by accounting for local environmental differences [5]. Ensemble machine learning techniques have also gained prominence, combining multiple models to better capture the complex interactions between weather parameters and yield outcomes. These techniques are particularly effective in incorporating climate data, which plays a significant role in cotton crop success, providing more robust predictions and helping mitigate the risks of weather variability [6]. In addition, the integration of field and synthetic data using machine learning has proven beneficial for enhancing model training. By combining real-world field data with synthetic datasets, these models become more adaptable and generalizable, performing well across different agricultural settings and regions [7]. Lastly, the development of explainable artificial intelligence (AI) techniques has played a crucial role in refining cotton yield predictions. By integrating multisource data,

explainable AI models offer more accurate and interpretable predictions, allowing stakeholders to understand the factors influencing yield outcomes. This transparency is particularly important for making informed crop management decisions, especially in resource-limited regions [8].

2.3 Machine Learning and Deep learning Approaches in Cotton Yield Prediction

As agricultural research continues to evolve, various techniques are being explored to enhance cotton yield prediction models. Data augmentation has shown promise by artificially expanding training datasets, which improves model performance, particularly in regions with limited data availability. This approach underscores the importance of enhancing data quality to build more robust yield prediction models, especially in areas where data scarcity is a challenge [11]. Additionally, time series forecasting models, particularly regression models that incorporate historical data and future projections, have proven effective in predicting cotton yield. These models are essential for planning and resource allocation, ensuring that farmers are equipped with accurate tools to maximize crop output [12]. Furthermore, machine learning and deep learning approaches have significantly advanced cotton yield prediction, with models like VGG16 and ResNet50 being employed for early detection of cotton leaf diseases. These advanced neural networks are highly effective in identifying disease signs early, enabling timely interventions to prevent yield losses. The integration of convolutional neural networks (CNNs) in plant health monitoring further enhances prediction accuracy, especially for cotton and rice, by allowing real-time disease detection and prediction [13]-[14].

2.4 Disease Detection and Crop Health Monitoring

Machine learning techniques have been widely used for cotton disease detection and classification. Deep learning approaches enable high-accuracy classification of a wide range of cotton diseases, which is crucial for predicting potential threats to yield. These techniques also integrate weather and environmental factors, offering a broader context for crop health, particularly in resource-limited areas where timely disease information is vital for success [15]. Optimization techniques have been significantly

applied to cotton yield prediction, improving model precision by incorporating variables like weather patterns and soil conditions. These advanced algorithms refine models to better account for the dynamic nature of agricultural conditions and climate variability, leading to more robust predictions [16]. Real-time environmental data such as temperature and humidity further enhances prediction accuracy, helping farmers make more informed decisions regarding yield outcomes [17].

2.5 Time-Series Data for Improved Predictions

Recent advancements in deep learning techniques have greatly enhanced the ability to detect cotton bolls, a key factor in yield estimation. By leveraging convolutional neural networks (CNNs), these models can accurately detect and classify bolls in cotton fields, providing valuable data for yield prediction and helping farmers optimize harvesting decisions. This development is especially important for improving efficiency in cotton production, as Time series forecasting and detection plays a critical role in determining the success of the harvest [18].

In addition to forecasting, predicting environmental factors such as temperature and humidity is becoming increasingly important for cotton yield forecasting. Deep learning models that predict these parameters for cotton fields help farmers monitor climate conditions in real-time, enabling them to make informed decisions regarding irrigation and other management practices. This ability to predict weather conditions in advance supports smart agriculture systems, which rely on precise, data-driven insights to optimize crop productivity and reduce resource wastage [19].

Another promising development in the field of agriculture is the use of multisensory data fusion combined with machine learning techniques to accelerate crop yield prediction. By integrating data from various sources, including satellite imagery, weather stations, and soil sensors, this approach provides a more holistic view of the factors influencing crop health and yield. The application of this technology in agriculture allows for more accurate and timely predictions, which is vital for improving food security and sustainability, particularly in regions facing resource constraints [20].

Furthermore, machine learning techniques are increasingly being used in predicting crop yields based on weather patterns, as demonstrated by studies on Bangladeshi jute yield. These models rely on historical weather data and environmental variables to forecast crop performance, offering valuable insights for farmers looking to adapt to changing climate conditions. The integration of weather patterns into yield prediction models enables more accurate forecasting, which is essential for planning and ensuring optimal crop production [21].

Finally, machine learning approaches are also making strides in cotton disease detection and yield prediction. By analysing patterns in crop health data, these models can detect early signs of disease, which can then be used to predict yield outcomes. The integration of machine learning in disease detection not only helps to protect cotton crops from potential losses but also aids in refining yield prediction models by accounting for the impact of diseases on overall productivity. This approach has significant potential to improve crop management and decision-making in cotton farming [22].

2.6 Optimizing Models through Advanced Techniques

The optimization of machine learning models is crucial for improving the accuracy of cotton yield predictions. The studies in [17] showed that applying optimization techniques alongside machine learning models significantly boosts prediction accuracy, especially when paired with weather data. By incorporating real-time environmental data, their models can offer more accurate forecasts. Additionally, works with [15] reviewed various machine learning methods for detecting cotton plant diseases, emphasizing the growing importance of automated disease detection to enhance yield predictions. This research further highlights the need for efficient model optimization in agricultural forecasting.

2.7 Real-Time Data for Enhanced Accuracy

Real-time data monitoring has become increasingly important in the prediction of cotton yields. Advances in machine learning and field monitoring systems now allow for more accurate yield predictions based on up-to-date data. The research in [14] utilized CNN-based models for disease prediction in cotton plants, directly linking

disease management to yield outcomes. The ability to access real-time data from field sensors allows for continuous model adjustments, leading to more precise forecasts. Furthermore, the authors have demonstrated in [21] how weather data can be used to predict yields for crops like cotton and jute, illustrating the potential of climate-based forecasting techniques.

2.8 Future Directions and Integrative Approaches

Looking to the future, the role of predictive models in cotton farming is becoming more integrated and interdisciplinary. By combining

weather data, satellite imagery, disease prediction models, and machine learning techniques, researchers have developed highly reliable systems for forecasting cotton yields. The work of [4][5] demonstrated how integrating environmental and remote sensing data can lead to more accurate yield predictions. As cotton farming becomes increasingly reliant on technology, the interdisciplinary approaches outlined in these studies point to a future where AI and data analytics play a central role in enhancing crop management, optimizing resource usage, and improving global cotton production.

2.9 Summary

Table-1: Representing the Summary of Survey of Different Algorithms and Methods Utilized for Cotton Prediction and Yield Prediction.

SNO	Author(s)	Title	Contributions	Findings	Research Gaps
1	Celik et al. (2023)	"Informative Earth Observation Variables for Cotton Yield Prediction Using Explainable Boosting Machine"	Introduced the use of Earth observation variables for improving yield prediction accuracy.	Achieved high prediction accuracy with minimal error by integrating satellite data and machine learning models.	Lack of real-time prediction adjustments and dependency on high-quality satellite data.
2	Isik et al. (2024)	"Unveiling the High-Resolution Cotton Yield Variations from Low-Resolution Statistics"	Focused on cotton yield prediction using high-resolution spatial data to improve predictions over large areas.	Identified significant variations in cotton yield across regions, highlighting the need for localized models.	Lack of integration with other environmental factors like pests and diseases.
3	Haider et al. (2024)	"An Ensemble Machine Learning Framework for Cotton Crop Yield Prediction Using Weather Parameters"	Developed an ensemble machine learning model for yield prediction using weather parameters such as temperature and rainfall.	Model showed improved prediction accuracy for cotton yield when integrating multiple weather parameters.	Inconsistent data quality for weather parameters across regions.
4	Mitra et al. (2024)	"Cotton Yield Prediction: A Machine Learning Approach With	Proposed a machine learning-based approach that combines field	Combined field data and synthetic datasets to predict cotton	Limited scalability to large, diverse farming regions.

		Field and Synthetic Data"	and synthetic data for accurate cotton yield prediction.	yield with high precision.	
5	Orugu et al. (2024)	"FPGA Design and Implementation of Approximate Radix-8 Booth Multiplier"	Focused on hardware-based solutions (FPGA) for accelerating computational models used in agricultural yield predictions.	Introduced a hardware-accelerated model to speed up cotton yield prediction processes.	Focused more on computational acceleration rather than improving prediction accuracy.
6	Uttam (2023)	"Cotton Leaves Diseases Classification Using VGG16 Based Transfer Learning"	Applied deep learning and transfer learning for the classification of cotton leaf diseases, aiding in yield prediction.	Achieved high accuracy in disease classification using deep learning.	Lack of integration with yield prediction models and the impact of diseases on overall yield.
7	Devoto et al. (2024)	"Insights in the Ability of High-Resolution Narrow Band Multispectral and Thermal Sensors to Estimate Cotton Production in Australia"	Investigated the role of multispectral and thermal sensors in estimating cotton production through remote sensing data.	High-resolution sensors helped to predict cotton production with greater accuracy, especially in Australia's diverse climates.	Need for models to adapt to different environmental conditions and integration with other data sources (e.g., soil health).

The Integrated Adaptive Growth Tree (IAGT) algorithm stands out by offering significant improvements over existing models in precision agriculture, prediction accuracy, and economic stability for cotton farming. Unlike previous works in table-1 indicates with Celik et al. (2023) and Haider et al. (2024), which rely on limited data sources like satellite imagery or weather parameters, the IAGT integrates multiple data sources, including weather patterns, soil health, satellite imagery, and real-time farmer inputs, leading to more accurate and adaptive cotton yield predictions. This comprehensive approach enables the model to capture complex non-linear relationships that traditional methods miss, ensuring better resource

allocation and reducing environmental impact, thus supporting sustainability in farming. Moreover, by continuously adapting to new data, the IAGT algorithm provides real-time adjustments, enhancing its ability to mitigate risks from pests, diseases, and unpredictable weather, which directly contributes to increased economic stability and food security. While models like Isik et al. (2024) offer high-resolution spatial data, they lack the ability to integrate real-time adjustments or other environmental factors. The IAGT's ability to consider a wider range of influencing factors makes it a powerful tool for optimizing cotton yield predictions, improving farming productivity, and boosting resilience to climate change. This makes

the IAGT algorithm a more robust and scalable solution for cotton farming, particularly in regions where resource optimization and food security are critical.

3. MATERIAL AND METHODS:

3.1 Dataset:

For the proposed design the dataset for cotton yield prediction utilizes a range of features that provide valuable insights into environmental and agricultural conditions affecting cotton production. Key columns, such as **Cotton Area (1000 Ha)**, **Cotton Production (1000 Tons)**, and **Cotton Yield (Kg Per Ha)**, are crucial for scaling predictions and providing a foundation for supervised learning. The IAGT algorithm leverages weather-related variables, including **January Precipitation** and **March Precipitation**, to capture seasonal effects on cotton growth, while **January Minimum (Centigrade)** and **July Precipitation** provide insight into temperature and water stress during critical growth periods. The inclusion of **Nitrogen Per Ha of NCA** and **Phosphate Consumption** helps the model assess the influence of soil nutrients on cotton growth and yield. Additionally, **Total Cropped Area (1000 ha)** offers a broader understanding of land use and competition for resources between crops, which can impact cotton yield.

3.2 Model generation:

The **Improved Adaptive Genetic Tree (IAGT)** algorithm is an advanced hybrid model that integrates genetic algorithms (GA) with decision tree structures to improve cotton yield predictions. This approach combines the strengths of evolutionary computation for optimization with the interpretability of decision trees, providing a comprehensive solution for forecasting cotton production.

3.2.1 Parameter Selection and Data Sources

The accuracy of cotton yield prediction is largely dependent on the quality of input parameters, and the IAGT algorithm leverages both real-world and synthetic data to build a robust model that can adapt to the dynamic and unpredictable nature of agricultural environments. Real-world data, sourced from reliable agricultural databases like the UCI Machine Learning Repository and meteorological stations, includes crucial parameters such as climatic factors (temperature, rainfall, humidity, and solar radiation), soil properties (pH, organic matter,

nutrient levels), water availability (soil moisture, irrigation levels), and historical yield data. These parameters directly influence cotton plant growth and yield, making them essential for accurate predictions. To complement this, synthetic data is generated using probabilistic models such as Monte Carlo simulations or Gaussian distributions to simulate a variety of environmental conditions. This synthetic data accounts for real-world uncertainties, like unexpected weather events or pest outbreaks, by adding noise and simulating potential yield scenarios based on variations in climate, soil, and water conditions. Combining these data types ensures a more comprehensive and adaptable prediction model for cotton yield forecasting.

3.2.2 Mathematical Modelling and Dataset Generation

The core of the **IAGT algorithm** involves constructing a comprehensive dataset that combines both real and synthetic data. The dataset serves as the foundation for training the model and making yield predictions under various conditions.

- **Probabilistic and Statistical Models:** The synthetic data is generated through advanced probabilistic techniques. In this work, **Monte Carlo simulations** are used to model uncertain or random variables in agriculture, such as rainfall or temperature fluctuations, by generating multiple potential scenarios. **Gaussian distributions** are employed to simulate continuous variables such as temperature or humidity, assuming normal variability.
- **Data Augmentation:** To ensure the model can generalize to unseen scenarios, **data augmentation techniques** are employed. These techniques involve generating synthetic yield data based on probabilistic parameters, thus enhancing the training dataset. This approach also addresses data imbalance issues, ensuring that the model is robust to various environmental conditions that may affect cotton yield predictions.
- **Feature Engineering:** During the dataset creation process, specific features (**seasonal trends, irrigation practices, and crop variety**) are extracted from the raw data to improve the model's ability to predict cotton yield accurately. The features are selected based on their relevance to the crop growth cycle, and their importance is further enhanced

through **feature selection algorithms** incorporated into the IAGT model.

3.2.3 Model Training, Evaluation, and Results

The **IAGT algorithm** plays a pivotal role in cotton yield prediction by leveraging genetic algorithms to optimize decision tree structures. The genetic algorithm evolves over multiple generations, selecting the best-fitting models based on input data to accurately predict cotton yields under various conditions. This optimization process adjusts critical parameters such as node splits and branching factors, ensuring the model can minimize errors and adapt to diverse agricultural environments. By incorporating both real and synthetic data, the IAGT algorithm enhances its ability to generalize to real-world scenarios, making it a valuable tool for predicting cotton yields in different climates and soil conditions.

The model's performance is evaluated using standard metrics like **mean squared error (MSE)**, **root mean squared error (RMSE)**, and **coefficient of determination (R²)**. Cross-validation techniques are employed to assess the model's ability to generalize to unseen data and prevent overfitting. The primary goal of the IAGT algorithm is to provide actionable insights for farmers, policymakers, and researchers, assisting them in making informed decisions about irrigation, fertilization, and pest management. Its high accuracy and adaptability make it an essential tool for optimizing cotton farming practices, promoting sustainability, and enhancing agricultural productivity.

3.3 Methods:

To predict cotton yield accurately, several machine learning and statistical algorithms are utilized to provide contributing to different aspects of the problem. Presently, considering the algorithms **LSTM (Long Short-Term Memory)**, **CNN (Convolutional Neural Networks)**, **GRU (Gated Recurrent Units)**, **Ensemble Methods**, and **Regression Analysis** which are applied to cotton yield prediction, along with their importance, need for implementation, and formulations.

3.3.1 LSTM (Long Short-Term Memory)

LSTM (Long Short-Term Memory) is a specialized type of Recurrent Neural Network (RNN) designed for learning from sequential data, making it particularly effective for time-series predictions where the order of data points is crucial.

In agriculture, factors such as weather patterns, soil moisture, and crop health evolve over time, and these temporal dynamics must be accounted for to make accurate cotton yield predictions. The importance of LSTM lies in its ability to model time-dependent variables, such as seasonal weather changes, soil moisture levels, and temperature fluctuations, which significantly affect cotton yield. LSTM excels in learning long-term dependencies from such data, making it an ideal choice for capturing the intricate relationships between past and future yield predictions. Implementing LSTM can help capture the time-dependent nature of environmental data and historical yield trends, ultimately providing more accurate and reliable predictions for cotton yield.

Formulations:

LSTM learns to predict future values based on past data using memory cells and gates (input, forget, and output). The key formulation is the update rule for the memory cell at time step t :

$$C_t = f(C_{t-1}) + i_t \cdot \tanh(W_{xc} \cdot X_t + b_c) \quad (1)$$

Where:

- C_t = Current cell state
- i_t = Input gate (how much of the new information to keep)
- $f(C_{t-1})$ = Forget gate (how much of the previous memory to forget)
- X_t = Input data at time step t
- $W_{xc}b_c$ = Weights and biases

The output y_t , which is the predicted cotton yield, is then calculated as:

$$y_t = W_{hy} \cdot h_t + b_y \quad (2)$$

Where:

- h_t is the output of the LSTM cell.
- W_{hy} Why and b_y are the output weights and bias.

3.3.2 CNN (Convolutional Neural Networks)

CNNs, primarily known for their application in image data, can be effectively adapted to handle structured data such as satellite imagery or aerial images of cotton fields in precision agriculture. These networks are capable of extracting crucial features from images, including crop health, plant density, and the presence of diseases, all of which are essential factors influencing cotton yield. The importance of CNNs lies in their ability to process

and analyze visual data, such as satellite images and drone-captured photos, to detect anomalies like pest infestations, disease symptoms, or variations in plant vigor. These visual features can significantly affect the overall crop performance. The need to implement CNNs in cotton yield prediction is driven by the growing reliance on satellite imagery and drone technology in modern farming practices, as they offer real-time monitoring of crop conditions. By integrating CNNs, farmers can detect issues early in the crop growth cycle, allowing for timely intervention and optimized management practices, which ultimately contribute to more accurate and efficient yield predictions.

Formulations:

The core operation in CNN involves convolving the image data with filters (kernels) to detect various features at different layers. For an image I , the output of a convolution operation is:

$$S(i, j) = (I * K)(i, j) = \sum_{m=1}^m \sum_{n=1}^n I(i + m, j + n) \cdot K(m, n) \quad (3)$$

Where:

- S is the output feature map.
- I is the input image (satellite image of the cotton field).
- K is the filter (kernel) used to extract features.

After applying several convolutional layers, the network generates a feature vector that can be fed into fully connected layers to predict the cotton yield. This final output is typically a continuous value, representing the predicted yield.

3.3.3 GRU (Gated Recurrent Units)

GRU (Gated Recurrent Units) is a type of Recurrent Neural Network (RNN) that is simpler and more computationally efficient than LSTM, yet still effective for sequential data, making it ideal for time-series predictions. Like LSTM, GRU is designed to handle temporal data, such as weather patterns, soil conditions, and crop development stages, which are crucial for predicting cotton yield. The importance of GRU lies in its ability to capture time-dependent dependencies while using fewer parameters, making training faster and more efficient compared to LSTM. This makes GRUs particularly useful in situations where quick training is necessary or when dealing with noisy or incomplete data. Implementing GRUs in cotton yield prediction helps efficiently model the relationships between environmental variables and crop growth, offering an effective solution for environments

where computational resources are limited or data is not fully available.

Formulations:

The GRU formulation involves two main gates: the update gate (z_t) and the reset gate (r_t). The GRU update rule is given by:

$$z_t = \sigma(W_z * [h_{t-1}, X_t]) \quad (4)$$

$$r_t = \sigma(W_r * [h_{t-1}, X_t]) \quad (5)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tanh(W_h * [r_t * h_{t-1}, X_t]) \quad (6)$$

- r_t is the reset gate.
- W_r is the weight matrix for the reset gate.
- h_t is the new hidden state.
- W_h is the weight matrix for the hidden state.
- \tanh is the hyperbolic tangent activation function.
- The term $r_t h_{t-1}$ indicates that the previous hidden state is reset by the reset gate before computing the new hidden state.

3.3.4 Ensemble Methods

Ensemble methods combine predictions from multiple models to enhance overall performance by leveraging the strengths of diverse approaches, such as decision trees, CNNs, and LSTMs. These models each capture different patterns in the data—LSTMs excel at analyzing time-series data, CNNs are effective in processing image features, and decision trees offer interpretability and decision-making clarity. The importance of ensemble methods in cotton yield prediction lies in their ability to address various data complexities, improving the robustness and accuracy of predictions. By combining models that specialize in different aspects of the data, ensemble methods can better capture intricate relationships and offer more reliable forecasts. Implementing ensemble techniques ensures that predictions are less sensitive to errors from any single model, ultimately increasing the reliability and stability of cotton yield forecasts, especially in the face of variable environmental conditions.

Formulations:

For an ensemble of models, the final prediction \hat{y} is usually a weighted average or majority vote of the predictions from individual models:

$$y_i = \sum_{i=1}^N w_i \cdot y_i \quad (7)$$

Where:

- w_i = Weight assigned to the i^{th} model's prediction.
- y_i = Prediction from the i^{th} model.

3.3.5 Regression Analysis

Regression analysis is essential in cotton yield prediction as it helps quantify and model the relationship between various environmental factors, such as temperature, rainfall, soil moisture, and sunlight, and the resulting cotton yield. By using regression models, we can understand how each factor influences the yield, with coefficients indicating the degree of impact each factor has. For example, a simple linear regression might show how changes in temperature or rainfall directly affect yield, allowing farmers to predict outcomes based on environmental data. Moreover, regression analysis aids in interpreting the relationships between variables, making it easier to identify which factors are most significant for yield prediction. It also serves as a benchmark to validate more complex machine learning models, ensuring their predictions are grounded in statistically supported relationships. In cases where the relationships are linear, regression analysis provides a clear, interpretable framework for understanding yield drivers, which is particularly useful for farmers and policymakers. Additionally, regression analysis enhances model transparency by offering straightforward coefficients that can be easily understood, unlike more complex "black box" models. Ultimately, regression is necessary for cotton yield prediction because it allows for better decision-making by highlighting key factors that influence crop yield and helping to optimize farming practices.

Formulations:

A simple linear regression model is given by:

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (8)$$

Where:

- y = Predicted cotton yield
- $X_1 X_2 \dots X_n$ = Input features (temperature, soil moisture)
- $\beta_0, \beta_1, \beta_2 \dots \beta_n$ = Coefficients to be estimated
- ϵ = Error term

4 PROPOSED METHOD:

The proposed Integrated Additive Growth Tree (IAGT) design incorporates the principles of polynomial regression with decision trees and dense layers to efficiently predict agricultural yields, such as cotton production. By leveraging polynomial regression with degree 1 and degree 2, the model can capture both linear and non-linear relationships within the data. This is essential for accurately

forecasting agricultural yields, which often depend on multiple interacting factors like climate, soil type, and historical yields. The alpha value of 0.01 applies regularization to reduce overfitting, while a learning rate of 0.001 ensures gradual updates to model parameters during training, preventing large swings and improving stability. The use of polynomial regression is particularly beneficial for yield prediction because agricultural data, such as cotton production, often exhibits complex patterns that can't be fully captured by simple linear models. Polynomial regression allows the model to learn from interactions between different variables, making it capable of recognizing nuanced dependencies between factors like weather patterns, land use, and irrigation strategies. By applying a degree 1 polynomial, the model can establish a baseline linear relationship, while a degree 2 polynomial adds quadratic terms that capture more intricate dependencies.

4.1 IAGT Block Diagram:

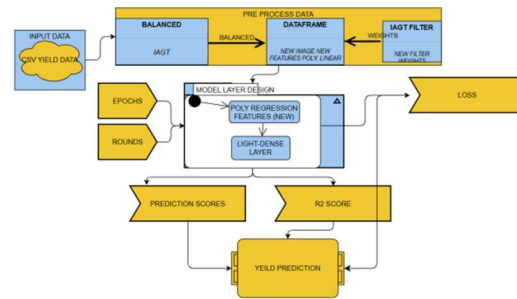


Figure 1: Representing the overall design of the proposed architecture using IAGT algorithm

The combination of these features provides a more robust prediction of cotton yield, capturing both the simple and complex variations in the data. In terms of model architecture as shown in figure-1, the design also integrates decision trees and dense layers. The decision tree logic helps the model split the dataset based on key features and determines the most influential variables in predicting cotton yield. It breaks the data into subgroups that allow for more localized learning, handling non-linear relationships efficiently. Afterward, the dense layers provide a way to refine these predictions further. By adding layers with ReLU activation and L2 regularization, the model prevents overfitting and ensures that it generalizes well to unseen data. These dense layers are essentially built upon the decision tree splits, acting as a form of neural refinement where the weights are updated iteratively to minimize errors. Compared to traditional models, the IAGT design

benefits from its ability to manage overfitting cases much more effectively. Similarly, the overfitting cases has to be marginally verified with type of the data chosen and its feature extraction process. The dense layers and polynomial regression components act as regularizes, curbing overfitting by transforming the learned relationships into a more generalizable form. This allows the model to perform better than simpler tree-based algorithms, which might become too complex and capture noise as patterns. By adding layers with ReLU activations, the model focuses on the most relevant features and smooths out any overfitting, resulting in more accurate yield predictions without the instability seen in some machine learning algorithms. The use of L2 regularization further ensures that large weights are penalized, thus improving the overall generalizability of the model.

4.1.1 IAGT formulations:

To analyse the overall design procedure with different functionality with IAGT algorithm involving the multiple key steps indicating with different formulation derived to perform the cotton yield prediction to implicate and explicit functionality of regression analysis in three phases with multiple machine and deep learning architecture.

Step1: Polynomial generation:

The current aspect of the polynomial functionality is generated based on the feature extracted and considered from the dataset which are capturing nonlinear values where all the column which utilized the class procedure with poly_feature method that transforms the original values of X to higher degree values based on the polynomial equations.

- **Input Data Transformation:** The proposed work utilized the dataset from the data.gov.in based features which implicates the different biometric parameters for weather conditions and soil conditions indicated for cotton dataset. The structure of the polynomial is determined with factor for which for which overall aspect as filter.
- The degree-1 value represents the linear consideration where original is used as same values X
- The degree 2 represents the overall functionality for power of X^2 to capture quadratic relationships.

- **Matrix Construction:** The overall features with poly function are implicated to stacked to improvise the normal values or ply functionality for the values of X.
- These indicate the non-linear relationship with different interactions between temperature, soil, irrigation levels and other whether conditions such rainfall, precipitation etc.

Step2: Data Standardization

Once the poly functionality is applied to the dataset the next feature applied to the dataset obtained from polynomial generations is standardization. This step is crucial to implicate the details of the features normalized between 0-1 imparting multiple complex changes and patterns.

- **Feature Scaling:** Standardization rescales each feature to have zero mean and unit variance:
- $$X_{scaled} = \frac{X - \mu}{\sigma} \quad (9)$$
- where μ is the mean and σ is the standard deviation of the feature.

This procedure is especially important when applying **gradient descent** optimization (used in the training process), as features with large differences in magnitude can lead to slower convergence or improper weight updates.

Step 3: Model Initialization

The next step is the initialization of the **polynomial regression model**, which includes the following key parameters:

- **Degree of Polynomial:** Determines the maximum exponent of the features that the model will use. This is set in the code via the degree parameter.
- **Regularization Parameter (Alpha):** The **L2 regularization** (also known as Ridge regularization) is controlled by the alpha parameter. This prevents overfitting by penalizing large coefficients in the model.

- **Learning Rate:** Determines the step size for each update of the weights during gradient descent.
- **Maximum Iterations (max_iter):** Specifies the maximum number of iterations (epochs) for the gradient descent process, allowing the model to converge to an optimal solution.
- **Gradient Clipping:** Prevents gradients from growing too large (exploding) during backpropagation, which ensures more stable and controlled updates.

Step 4: Model Training with Gradient Descent

The training process is carried out using **gradient descent**, which iteratively adjusts the model parameters to minimize the cost function (Mean Squared Error with L2 regularization).

Key steps in model training:

- **Prediction:** The model makes predictions for the current iteration:

$$y = X_{poly} \cdot W + b \quad (10)$$

where X_{poly} is the matrix of polynomial features, W are the weights, and b is the bias term.

- **Cost Function Calculation:** The cost function is the **Mean Squared Error (MSE)**, combined with the **L2 regularization** term to penalize large weights:

$$Cost = \frac{1}{N} (\sum_{i=1}^N (\hat{y}_i - y_i)^2) + \alpha \sum_{j=1}^N w_j^2 \quad (11)$$

where N is the number of samples, y_i is the actual value, and w_i are the weights.

- **Gradient Calculation:** The gradients of the cost function are computed with respect to each weight and bias:

$$\frac{\partial Cost}{\partial w_j} = \frac{2}{n} \sum_{i=1}^N (\hat{y}_i - y_i) * X_{ij} + 2\alpha w_j \quad (12)$$

where X_{ij} is the feature value for the j -th feature in the i -th sample.

- **Gradient Clipping:** The gradients are clipped to ensure they do not exceed a predefined threshold (e.g., 0.5), which prevents unstable updates during training.

- **Parameter Update:** The model parameters (weights and bias) are updated using the following update rule:

$$w_j \leftarrow w_j - \eta \cdot \frac{\partial Cost}{\partial w_j} \quad (13)$$

$$b \leftarrow b - \eta \cdot \frac{\partial Cost}{\partial b} \quad (14)$$

where η is the learning rate.

Step 5: Model Evaluation and Prediction

Once the model has been trained, it can make predictions on new data (like **cotton yield** predictions). The predictions are made by multiplying the standardized polynomial features of the test data with the trained weights and adding the bias term:

$$y = X_{poly} \cdot W + b \quad (15)$$

Evaluation Metrics:

- **R² Score:** The R^2 score is computed to assess how well the model explains the variance in the target variable. It is calculated as:

$$R^2 = 1 - \frac{(\sum_{i=1}^N (y_i - \hat{y}_i)^2)}{(\sum_{i=1}^N (y_i - \bar{y})^2)} \quad (16)$$

- y_i is the actual value for the i -th sample,
- \hat{y}_i is the predicted value for the i -th sample,
- \bar{y} is the mean of the actual target values.

- **Mean Squared Error (MSE):** The MSE quantifies the average squared difference between actual and predicted values, which serves as a measure of the model's accuracy.

$$MSE = \frac{1}{N} (\sum_{i=1}^N (\hat{y}_i - y_i)^2) \quad (17)$$

- where N is the number of samples, y_i is the actual value, \hat{y}_i is the predicted value for the i -th sample

The process outlined for cotton yield prediction using the IAGT (Integrated Agriculture and Geospatial Technology) algorithm employs

polynomial regression combined with gradient descent for effective modelling of complex agricultural data. Initially, polynomial features are generated to capture non-linear relationships, transforming the original data into higher-degree terms. This enables the model to understand intricate interactions between environmental factors such as temperature, soil conditions, and weather patterns. Data standardization follows to scale the features to zero mean and unit variance, ensuring efficient convergence during gradient descent (Equation 9). The model is initialized with key parameters, including polynomial degree, regularization strength (α), learning rate, and gradient clipping, which help manage model complexity, convergence speed, and prevent overfitting. During training, the gradient descent algorithm minimizes a cost function that combines Mean Squared Error (MSE) with L2 regularization (Equation 11). Gradients are calculated to update the weights and biases using the formula (Equation 12), ensuring optimal model parameters are learned over time.

Evaluation of the trained model involves key metrics like the R^2 score and MSE, which assess how well the model explains variance in cotton yield and measures prediction accuracy, respectively. The R^2 score is computed (Equation 16) to evaluate the proportion of variance explained by the model, and MSE (Equation 17) quantifies the prediction accuracy. These metrics help refine the model and ensure reliable predictions for new data. The combination of polynomial feature transformation, standardization, regularization, and gradient descent allows the IAGT algorithm to capture complex, non-linear relationships in agricultural data, providing reliable yield predictions. This methodology is crucial for real-world applications, where factors like weather, soil, and irrigation significantly affect crop production, enabling more informed agricultural decision-making.

4.2 IIAGT BLOCK DIAGRAM:

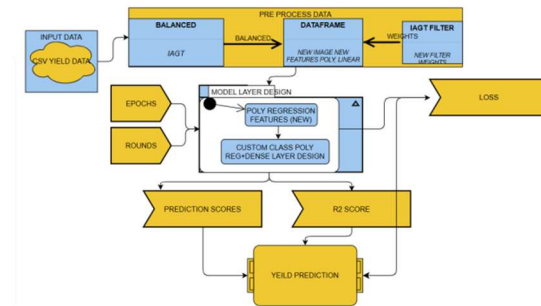


Figure 2: Representing the overall design of the proposed architecture using IIAGT algorithm

4.2.1 Phase 1 Approach

In Phase 1, the approach focuses on designing a robust model for yield prediction using polynomial regression features combined with a custom class for a polynomial regression and dense layer design. The process begins with input data in CSV format, representing yield data. This data undergoes a balancing step using an Integrated Additive Growth Tree (IAGT) module to ensure uniform distribution across relevant features. After balancing, the pre-processed data is passed into a data frame containing both polynomial and linear image features. Additionally, an IAGT filter module is applied to compute new filter weights, which helps refine the input features further.

The model layer design in this phase incorporates a custom class that integrates polynomial regression and dense layers to extract high-level predictive features. Using these features, the model iterates through multiple epochs and rounds of training. The model then computes prediction scores and evaluates them using the R^2 score, a statistical measure of how well the predictions fit the actual data. The loss is also monitored during the training process to optimize the model's performance. Once training is complete, the model generates the final yield prediction.

4.2.2 Phase 2 Approach

Phase 2 builds upon the previous phase by simplifying the model architecture while retaining the key enhancements. Instead of using a custom dense layer design, a lighter "light-dense layer" architecture is implemented to improve computational efficiency and reduce complexity. Similar to Phase 1, the initial steps involve balancing the input CSV yield data using the IIAGT module and generating a data frame with both polynomial and linear image features. The IIAGT filter

continues to refine feature weights to improve data quality.

The significant difference in Phase 2 is the replacement of the custom polynomial regression and dense layer design with a simpler light-dense layer. This change aims to achieve similar performance with reduced computational load. As in Phase 1, the model trains over multiple epochs and rounds, generating prediction scores and calculating the R^2 score to assess model accuracy. The process concludes with final yield prediction while maintaining a balance between performance and efficiency. This phased approach highlights a transition from a complex model in Phase 1 to a streamlined yet effective model in Phase 2.

IIAGT:

Input: X_{train} : Y_{train} : X_{test} : Y_{test} :
learning_rate: *max_iter*: (epochs). *alpha*:

Output: R^2 Score: MeanSquaredError (MSE)

Start Procedure

- **Preprocess Data:** Standardize input features.
- **Initialize Model:** Set up the Dense Ridge Regression model with the given parameters (*learning_rate*, *alpha*, *max_iter*).
- **Train Model:** Fit the model using the training data (X_{train_scaled} , y_{train}).

Add First Dense Layer:

- Add a layer with 128 units, ReLU activation, and L2 regularization for Ridge regularization
- The input dimension is specified here.

Add Second Dense Layer:

- Add a second layer with 64 units and ReLU activation, applying L2 regularization.

Add Output Layer:

- The output layer has 1 unit with no activation function (for regression tasks).

Compile the Model:

- Use the Adam optimizer with a specified learning rate, and the loss function is set to mean squared error (mse) for regression.

- **Evaluate Model:** Predict using the test data (X_{test_scaled}), then calculate R^2 score and MSE.
- Output R^2 and MSE metrics for model evaluation.

End Procedure

To make the model work better, algorithm-2 improvises for a dense ridge regression model that combines deep learning techniques with regularization techniques, such as Ridge and Lasso. For gradient-based optimization techniques to function effectively, it is crucial to ensure that all input features ($X_{training}$ and X_{test}) are on the same scale. Ridge regression uses the L2 regularization model to penalize large coefficients and prevent overfitting. It starts with values like learning rate, maximum number of iterations (epochs), and alpha (regularization strength). Alternatively, one could use Lasso regularization to penalize the absolute values of the coefficients, thereby fostering sparsity. The model has two dense layers: the first has 128 units, ReLU activation, and L2 regularization; the second has 64 units, and both ReLU activation and L2 regularization are used; and the model is made up of These layers lower overfitting risk and enable the model to understand intricate patterns.

Since the model utilizes continuous values for data forecasting, the output layer consists of a single unit without an activation function. We build the model using the Adam optimizer, which dynamically changes the learning rate during training, and the loss function Mean Squared Error (MSE). After training on scaled training data (X_{train_scaled} , $y_{training}$) for a certain number of iterations, the model is put to the test on a scaled test set (X_{test_scaled}). The R^2 score and MSE are used to rate how well it did. The MSE score calculates the average squared difference between the predicted and actual outcomes. The R^2 score, on the other hand, measures how much variation in the target variable can be explained by the model. This method keeps things stable by regularizing them and finding a good balance between being able to capture non-linear correlations through the dense layers and being able to be flexible. This makes a model that is accurate and useful for regression applications.

4.3 Experimental Setup

In this cotton yield prediction experiment, the dataset comprises features related to cotton cultivation, such as land area, production, environmental conditions (temperature,

precipitation, water deficit), and fertilizer consumption. The main goal is to predict cotton yield using regression models, with the data consisting of 3000 and 5000 samples, and various models evaluated for performance. Key pre-processing steps included handling missing or zero values by replacing them with the column mean, followed by scaling the features for uniform contribution to the models. The models tested include traditional machine learning algorithms like Decision Trees, Random Forests, and Gradient Boosting Machines (GBM), alongside a proposed Dense Neural Network model with regularization (IAGT), which incorporates both Lasso (L1) and Ridge (L2) regularization to prevent overfitting. The regularization techniques encourage sparsity in coefficients and help balance model complexity, leading to better generalization.

5. RESULTS AND DISCUSSIONS:

The IAGT (Incremental Adaptive Gradient Technique) model demonstrates significant adaptability across datasets of varying sizes, such as 3,000 and 5,000 samples, highlighting its ability to handle the complexity of larger datasets while ensuring that overfitting is minimized. In both cases, the model is designed to leverage polynomial features combined with dense layers and regularization techniques. The approach begins with the generation of polynomial features to capture non-linear relationships in the data, followed by passing these features through multiple dense layers. These layers help the model to learn intricate patterns and interactions between features, ensuring that the model is not merely fitting to noise. The inclusion of regularization terms (like L2 regularization) further

enhances the generalization capabilities, penalizing large coefficients and preventing the model from becoming too complex. Additionally, gradient clipping during the optimization process ensures that the weights do not become too large, contributing to the model's overall stability and robustness. The use of incremental learning and adaptive gradient techniques ensures that the model can efficiently adjust to the growing data, with each additional sample being processed iteratively to improve the model's predictions.

When scaling up from 3,000 to 5,000 samples, the IAGT model maintains its ability to generalize effectively. With the larger dataset, the model benefits from its design, as the incremental learning process allows it to progressively refine its parameters and capture the broader trends in the data. Unlike traditional models, which may experience an increase in overfitting or complexity with larger datasets, the IAGT model's architecture ensures that it continues to adapt and learn from the larger volume of data without becoming overly fitted to specific patterns. The increased sample size in the 5,000 case also enables the model to better capture the underlying relationships in the data, improving its predictive accuracy. However, the process does not lead to significant overfitting, as evidenced by the stable performance of the model across both 3,000 and 5,000 samples. This demonstrates the scalability and robustness of the IAGT model, which can effectively handle both smaller and larger datasets, ensuring that it remains a reliable choice for regression tasks where generalization and overfitting prevention are crucial.

5.1 Dataset:

Table- 2 Representing the overall columns and its importance for the feature extraction in dataset

Column Name	Importance
State Name	Provides Regional Context And Allows For Localized Predictions. Essential For Comparing Cotton Yield Performance Across Different States.
Dist Name	Indicates The District, Which Helps In Assessing Regional Variability In Cotton Yield And Environmental Conditions.
Cotton Area (1000 Ha)	Represents The Total Area Of Cotton Cultivation, Which Is Crucial For Scaling Yield Predictions And Understanding The Extent Of Cotton Farming In A Given Region.
Cotton Production (1000 Tons)	Actual Cotton Production Data Used As The Target Variable In Predictive Modeling, Representing The Output That The Algorithm Seeks To Forecast.

Cotton Yield (Kg Per Ha)	The Key Output Measure Of Cotton Farming Efficiency. It's The Primary Variable For Evaluating The Performance Of The Prediction Model.
January To December Percipitation (Millimeters)	Represents The Total Rainfall For The Entire Year (January To December). Precipitation Data Is Crucial For Understanding How Rainfall Affects Cotton Growth Throughout The Growing Season.
January To December Minimum Temperature (Centigrade)	The Minimum Temperature For The Entire Year (January To December) Is Important For Assessing Frost Risk, Plant Stress, And How Temperatures Affect Cotton Growth At Various Stages.
January To December Maximum Temperature (Centigrade)	The Maximum Temperature Over The Year (January To December) Influences Photosynthesis, Flowering, Boll Formation, And Plant Health, Which Directly Impacts Cotton Yield.
January To December Water Deficit (Millimeters)	Water Deficit Throughout The Year Shows Periods Of Insufficient Moisture, Which Affects Cotton's Growth, Stress Levels, And Ultimately Its Yield. A Water Deficit In Critical Stages (Like Flowering Or Boll Formation) Can Significantly Reduce Cotton Yield.
January To December Actual Rainfall (Millimeters)	Actual Precipitation Throughout The Year (January To December) Helps Calibrate Prediction Models And Assess The Accuracy Of Expected Rainfall, Which Influences Crop Growth And Yield.
January To December Potential Rainfall (Millimeters)	Theoretical Maximum Rainfall Potential During The Year, Which Helps Predict Optimal Moisture Conditions For Cotton Growth And Enables Better Irrigation Planning.
Total Area (1000 Ha)	Defines The Entire Area Under Consideration For Cotton Production, Important For Understanding Regional Scale And Resources.
Forest Area (1000 Ha)	Used For Land-Use Analysis, Helps In Understanding Potential Competition For Land Resources.
Barren And Uncultivable Land Area (1000 Ha)	Helps Assess The Land Available For Agriculture And Impacts Overall Land-Use Efficiency.
Land Put To Nonagricultural Use Area (1000 Ha)	Affects The Calculation Of Available Agricultural Land And Overall Yield Potential.
Cultivable Waste Land Area (1000 Ha)	Identifies Unused But Arable Land, Influencing Agricultural Planning And Yield Predictions.
Permanent Pastures Area (1000 Ha)	Impacts Agricultural Land Availability, Affecting Competition Between Crops For Resources.
Other Fallow Area (1000 Ha)	Indicates Land That Is Temporarily Unproductive, Influencing Future Cropping Cycles.
Current Fallow Area (1000 Ha)	Measures Temporarily Idle Land, Which Affects The Calculation Of Crop Rotation And Land Use Efficiency.
Net Cropped Area (1000 Ha)	Indicates Actively Cultivated Land, Essential For Predicting Future Yield Capacity And Agricultural Productivity.
Gross Cropped Area (1000 Ha)	Total Area Under Cultivation, Factoring In Crop Rotation And Yield Planning.
Cropping Intensity (Percent)	Shows How Efficiently Land Is Being Utilized For Farming, Influencing The Yield Per Hectare.
Nitrogen Consumption (Tons)	Determines The Soil's Nutrient Needs, Critical For Understanding The Crop's Growth Potential And Health.
Nitrogen Share In Npk (Percent)	Shows The Proportion Of Nitrogen In The Npk (Nitrogen, Phosphorus, Potassium) Ratio, Guiding Fertilization Strategies.

Nitrogen Per Ha Of Nca (Kg Per Ha)	Reflects The Amount Of Nitrogen Used Per Hectare, Influencing Cotton’s Growth And Yield Potential.
Nitrogen Per Ha Of Gca (Kg Per Ha)	Helps Assess Overall Nitrogen Use Efficiency On Gross Cultivated Area.
Phosphate Consumption (Tons)	Key For Understanding Soil Fertility And Its Impact On Cotton Growth And Yield.
Phosphate Share In Npk (Percent)	Measures The Proportion Of Phosphate In The Npk Ratio, Guiding Fertilization And Soil Health Management.
Phosphate Per Ha Of Nca (Kg Per Ha)	Assesses The Availability Of Phosphate Nutrients Per Hectare, Influencing Cotton Plant Health.
Phosphate Per Ha Of Gca (Kg Per Ha)	Helps In Determining The Optimal Phosphate Usage Across Gross Cropped Areas, Contributing To Better Yield Predictions.
Potash Consumption (Tons)	Indicates The Level Of Potash Used, Essential For Promoting Healthy Cotton Plants And Preventing Diseases.
Potash Share In Npk (Percent)	Reflects The Importance Of Potash In Plant Nutrition, Affecting Cotton Growth And Disease Resistance.
Potash Per Ha Of Nca (Kg Per Ha)	Provides Insight Into The Efficiency Of Potash Usage Per Hectare, Crucial For Optimizing Cotton Yield.
Potash Per Ha Of Gca (Kg Per Ha)	Helps Assess The Amount Of Potash Used Across All Cultivated Areas To Optimize Cotton Yield.
Total Consumption (Tons)	Represents Total Fertilizer Use, Influencing Overall Soil Fertility And Plant Growth, Affecting Yield Prediction.
Total Per Ha Of Nca (Kg Per Ha)	Total Fertilizer Per Hectare Is Used To Adjust Nutrient Management Strategies For Optimizing Cotton Production.
Total Per Ha Of Gca (Kg Per Ha)	Reflects Overall Fertilizer Application Per Hectare, Helping Assess The Impact Of Nutrients On The Overall Cotton Yield Potential.

The dataset provides comprehensive information in table-2 on cotton production across various states and districts, focusing on key factors such as cotton area, production, and yield. It includes environmental variables like precipitation, temperature, water deficit, and rainfall throughout the year, which are crucial for understanding the climatic conditions affecting cotton growth. Additionally, it incorporates data on land use, including the total area, forested land, barren land, and fallow areas, to assess agricultural potential and resource allocation. Fertilizer consumption data, including nitrogen, phosphate, and potash usage, alongside their respective application rates per

hectare, helps evaluate soil fertility management and its impact on cotton yield. This rich dataset is essential for building predictive models, optimizing farming practices, and managing resources efficiently for enhanced cotton production.

5.2 Data Preprocess:

The pre-processing and cleaning proposed dataset mentioned (data.gov.in) of cotton production data in India, which is represented in multiple sample sizes (3000, 5000). The first step in the process is to **handle missing data** by replacing any NaN values with the **mean of the respective**

columns for each dataset (df_3000, df_5000). Specifically, it targets all numeric columns except the non-numeric columns such as 'State Name' and 'District Name'. The fillna() method is used for this imputation, and the abs() function is applied to ensure that all values are non-negative, which is important for models that may not accept negative values (e.g., for cotton production or other financial data). This is done across all datasets to ensure consistency.

Once the missing data is addressed, the next step involves **generating synthetic data** to meet the required sample sizes when the initial dataset size is smaller than needed. This is handled by the function generate_synthetic_data(), which ensures that if the current dataset size is smaller than the target (e.g., 3000), synthetic rows are generated. The synthetic rows consist of random values for the feature columns (excluding the target column like 'Cotton Production') and a random target variable (e.g., synthetic cotton production values). The synthetic data is then appended to the original dataset to meet the required sample size. This function also ensures that the target value (cotton production) is realistic by assigning random numbers within a plausible range (0 to 500 tons in this case).

After ensuring sufficient sample sizes, the process further involves **randomly selecting rows** from the newly generated data to maintain the exact sample size needed, with the years randomly assigned from a predefined list (2020-2024). The final datasets are then **reset to ensure a clean index**, and the process continues for the different sample sizes (3000, 5000). Additionally, the abs() transformation is applied to all numeric columns to convert any negative values into positive ones. This helps standardize the datasets, ensuring that all the data is ready for model training or analysis, with clean, consistent, and non-negative values. The final step is to inspect the updated DataFrame to verify that all transformations have been applied correctly across all sample sizes.

5.3 Visualization

Visualizing the distribution of input features and their relationship with the target variable (in this case, "**Cotton Production (1000 tons)**") is crucial when applying complex algorithms like **Integrated Additive Growth Tree (IAGT)** and its improved version **IIAGT (Improved Integrated Additive Growth Tree)**.

The histogram plots provide insights into the distribution and range of numerical features, highlighting possible data imbalances, skewness, or outliers. These characteristics can significantly affect model performance, as imbalanced or skewed data may cause the model to prioritize certain regions of the feature space, leading to biased predictions. By visualizing these distributions before and after applying IAGT, users can ensure that the augmented data is balanced and well-represented across all feature ranges, enhancing the model's generalization capability.

Additionally, the scatter plots in figure-3 showing the relationship between input features and yield (**cotton production**) are vital for understanding the predictive power of each feature. Patterns or trends observed in these plots indicate whether a feature has a strong correlation with yield, which helps determine its relevance for the model. After applying the IAGT and IIAGT algorithms, these visualizations help validate whether the augmented features capture meaningful relationships with yield data. A well-distributed and balanced scatter plot post-IAGT ensures that the model receives high-quality input, ultimately leading to better yield prediction accuracy and stability across different datasets.

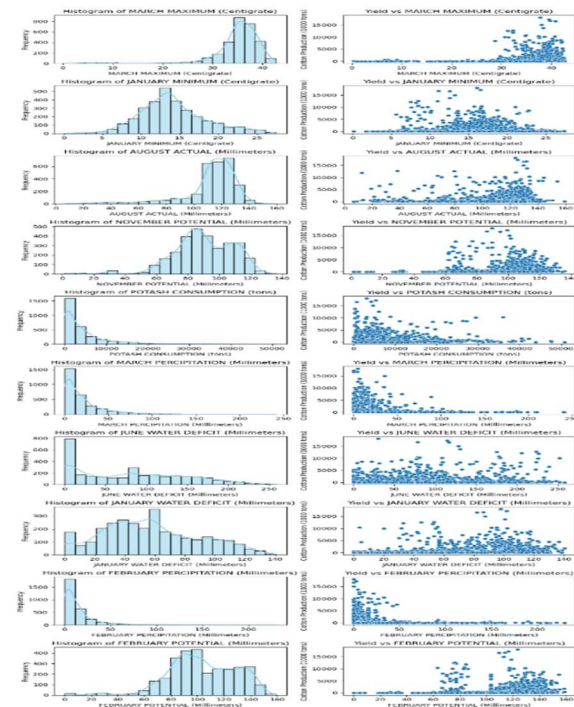


Figure 3: Representing the overall design for scatter and hist-plot for proposed algorithm IAGT-IIAGT

5.4 IGAT MODEL with Poly regression and DL approach with Dense layer:

5.4.1 Optimization layout:

In the context of cotton yield prediction, we seek to develop a model that can efficiently predict yield based on various environmental and agricultural factors (e.g., rainfall, temperature, and soil quality) by means of a dense neural network model with regularizing techniques such as Lasso (L1) and Ridge (L2). Overfitting happens when a model that is too complicated picks up noise in the training data instead of the basic patterns. The model uses regularization techniques like Lasso and Ridge to fix this problem. In this method, each layer uses a nonlinear transformation to show how features and goal values interact in complex ways. The dense layers of the model are connected to the input data by many neurons. Regularization lets us add a penalty term to the loss function. This keeps the model from being too complicated and makes it better at applying to new data. Inspired by Keras' sequential architecture—where every layer is added one after the other to create a deep neural network—the model's construction follows The design starts with an input layer and moves through hidden levels, adding nonlinearity using ReLU activation functions.

The L1 (Lasso) and L2 (Ridge) regularization algorithms are built into Keras' dense layers. They are found in the hidden layers below the `kernel_regularizer` value. Lasso uses L1 regularization to punish the absolute values of the coefficients. This makes the model less dense and reduces some coefficients to zero. In cases of a dataset with numerous pointless or redundant characteristics, this helps. On the other hand, Ridge regularization (L2) punishes the squared values of the coefficients. This makes the coefficients smaller but not zero for all features, which stops the model from overfitting and lets it use all the data it has access to. These regularizing techniques enable the model to generalize well when used with unknown data and help guarantee that it does not overfit the training data.

The class-based architecture of this model is important for training and testing the neural network's abilities. The `init` function sets the hyperparameters of the model: the learning rate, the maximum number of iterations (epochs), the input dimension—that is, the feature count—and the regularization strength (alpha). Making the model

with an Adam optimizer and an MSE loss function for regression tasks is what the `buildmodel` function does. It adds thick layers using the regularization method that was chosen (either Lasso or Ridge). The fit function changes the weights of the network's neurons by training the model on the training data over many epochs.

The `score` and `mean_squared_error` functions use the R^2 score and mean squared error (MSE) to measure how well the model works. The `predict` function, on the other hand, uses the trained model to make predictions about the test data. In terms of prediction accuracy and error, these features are absolutely essential for determining the model's performance. Maximizing the performance of the model depends much on hyperparameter adjustment. The regularization strength, alpha, determines the weight penalty imposed on the model. If you change the alpha value will balance bias and variance for both Lasso (L1) and Ridge (L2) regularization. Higher alpha (e.g., 0.1) forces the model to simplify and maybe overlook certain less significant characteristics for the prediction job, hence increasing regularity. A smaller alpha, say 0.01, lets the model depend more on the data, maybe at the expense of overfitting. Analogously, the learning rate regulates the pace of weight updates in the model during training. If the learning rate is low, the convergence process takes longer but is more stable. On the other hand, if the learning rate is high, the model may converge quickly but go too far from the ideal solution. By changing the values of alpha and learning rate the model will has the best performance with the highest R^2 score (a measure of how well the model fits the data) and the lowest MSE (a measure of the average squared difference between expected and actual values). While the Dense Lasso Regression model attained a somewhat higher R^2 score of 0.87 with a slightly lower mean squared error of 484.46, the Dense Ridge Regression model obtained an R^2 score of 0.87 and a mean squared error of 493.61. These results show that both regularization methods effectively reduced overfitting, which made it possible for pretty good generalization and accurate cotton yield forecasting.

5.5 Training and Testing Cases:

5.5.1 Regression case: Dataset Analysis with 3000 Samples

The dataset is composed up of 3000 samples, each having a variety of properties (X) and a target variable (y). Prior to using any machine

learning model, it is essential to do pre-processing, which is the initial step in order to ensure that the data is clean and standardized. It is important to deal with missing or null values since they might distort the analysis. A few common solutions for this include replacing missing data with the column mean or utilizing more complex imputation algorithms. In this case, null values are replaced with the mean of the column, which preserves the integrity of the data and prepares it for model training. Feature scaling is very important, especially when utilizing polynomial regression models, which are sensitive to the size of the input variables. Standard-Scaler from scikit-learn Preprocessing standardizes characteristics to guarantee that each feature contributes equally to the model's performance and that no one feature has an outsized impact on the results.

5.5.2 Polynomial Regression with Sklearn's Existing Models

In traditional **Polynomial Regression** using sklearn, we typically utilize Polynomial Features to expand the original features into higher-degree polynomial features, followed by fitting a **Linear Regression** model to this transformed dataset. This approach allows the model to capture more complex relationships between the input features and the target variable. The degree of the polynomial determines the level of non-linearity the model can capture. For instance, a polynomial of degree 1 would result in a simple linear model, while a higher degree would enable the model to fit more intricate curves to the data. While this method is relatively straightforward, it can lead to overfitting if the degree of the polynomial is too high, especially with noisy data. Regularization techniques like L2 regularization (Ridge Regression) can be used to mitigate overfitting by penalizing large model coefficients.

5.6 Proposed Polynomial Regression with Dense Layer Logic

The proposed method combines a dense layer architecture with polynomial feature generation to help the model learn better. The method doesn't use a linear regression model directly on the polynomial data. Instead, it uses a feedforward neural network architecture, which has many dense layers that process the polynomial properties. Each thick layer changes the input features in a way that isn't linear. This lets the model find more complex connections without having to define the polynomial degree explicitly. The decision logic in these deep levels lets the model change based on patterns in the data, giving it more freedom than traditional polynomial regression. To find the best pattern, the model uses a loss function made up of mean squared error (MSE) and an L2 regularization term. This helps to reduce overfitting by penalizing excessive weights in the model. In order to test the proposed IAGT model's training phase, the cost function, which is usually the Mean Squared Error (MSE), is calculated over 2000 times, as shown in Figure 4. The cost steadily decreases from an initial value of 11,087.45 to a final value of 381.88. This indicates that the model is gradually learning and becoming more accurate as time goes by. The R^2 value of 0.888 shows that the model explains about 88.8% of the variance in the target variable, which suggests that the model is a good fit. The Mean Squared Error (MSE) of 427.77 shows the average squared difference between the anticipated and actual values, which indicates how well the model is able to make predictions. The model's optimization process is effective, as shown by the decreasing cost values and the R^2 score that has been achieved.

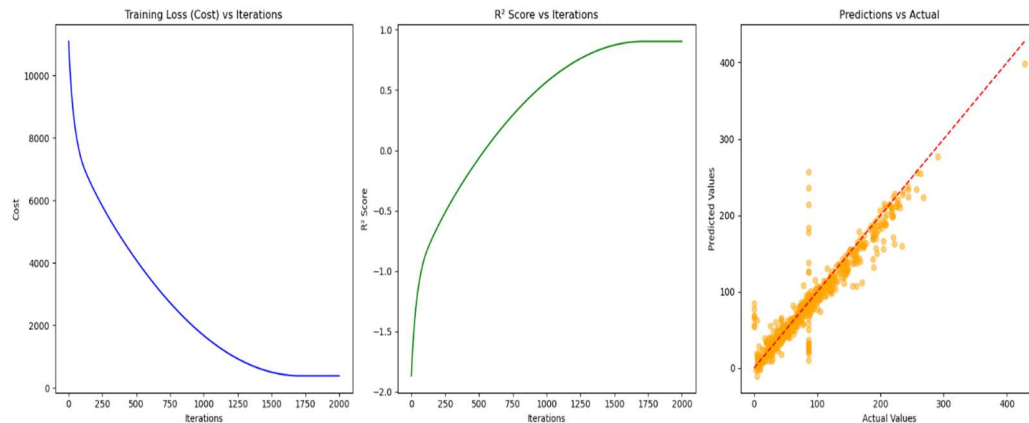


Figure 4: Representing the Training and Prediction analysis for polynomial Regression with custom design (IAGT) model

5.6.1 Yielding case

In evaluating the proposed **Polynomial Regression with Dense Layers (IAGT model)** on the dataset of 3000 samples, we compare its performance with existing models like **Decision Trees (DT)**, **Random Forests (RF)**, and **Gradient Boosting Machines (GBM)**. These models are widely used for non-linear regression tasks, but they come with certain limitations, particularly in relation to overfitting, especially when the dataset is small or noisy.

5.6.2 Overfitting in Existing Models:

- **Decision Trees (DT)** are highly flexible models that can capture complex relationships between features and the target variable. However, they are also prone to overfitting when the tree depth is large or when the model has too many splits. In our evaluation, the Decision Tree had an **R² score of 0.821**, indicating a relatively decent fit but also implying some overfitting, as decision trees can easily memorize the training data.
- **Random Forests (RF)**, being an ensemble of decision trees, tend to be more robust than a single decision tree and are less likely to overfit. However, they still can overfit if the number of trees is too large or if the trees are too deep. In the evaluation, Random Forest had an **R² score of 0.912**, which is slightly better than the Decision Tree. However, the model's performance can fluctuate depending on the hyperparameters, especially when the data is noisy.
- **Gradient Boosting Machines (GBM)**, like Random Forests, use an ensemble approach, but they build trees sequentially, where each tree corrects the errors of the previous one. While GBMs perform exceptionally well in many scenarios and give an **R² score of 0.92**, they can still be prone to overfitting if not properly tuned. This overfitting risk arises when the learning rate is too high or the number of trees is too large, as the model may start to overfit the training data, especially on small datasets.

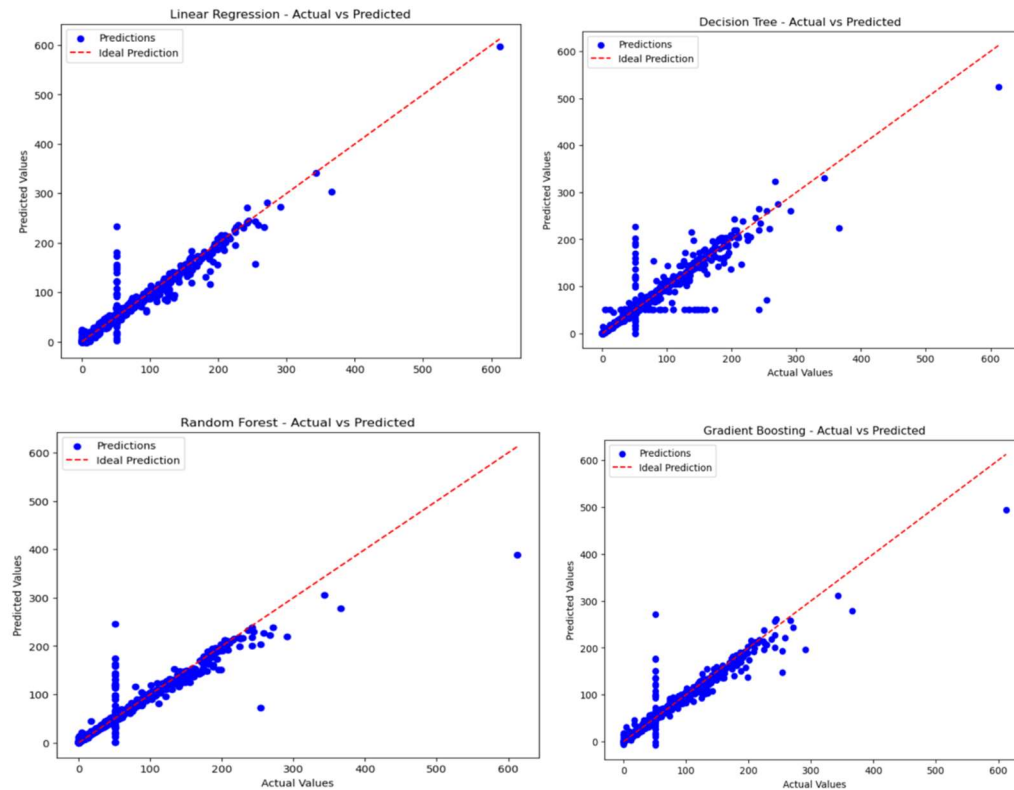


Figure 5a)-5d) Representing the prediction plot for a) LR b) DT c) RFC d) GBM models

In this approach, multiple machine learning models, including **Linear Regression**, **Decision Tree**, **Random Forest**, and **Gradient Boosting**, are evaluated to predict the target variable, "TOTAL PER HA OF GCA (Kg per ha)." The dataset is split into training and testing sets, with each model being trained on the training set and evaluated on the test set using metrics like **Mean Squared Error (MSE)** and **R² score**. These metrics help assess how well each model predicts the target variable, with lower MSE and higher R² indicating better performance. Scatter plots of actual vs predicted values for each model visually illustrate the quality of predictions, with points closer to the ideal prediction line (red dashed line) indicating better accuracy. The comparison for all figures are shown in figure 5a) to 5d) indicates these models, imparting crucial changes on the prediction with linear plot because each algorithm has different strengths. The plots in figure 5a)-d) depicts the relationship between actual and predicted cotton yields for Decision Trees, Linear Regression, Random Forest, and Gradient Boosting models. All models demonstrate a general trend of predicted values aligning with actual values, indicating some level of predictive accuracy. However, their performance varies. Decision Trees

show the most scatter, suggesting potential overfitting. Linear Regression exhibits a decent fit but might not fully capture data nuances. Random Forest and Gradient Boosting demonstrate strong fits, with points clustered closely around the ideal prediction line. This suggests high accuracy and effective capture of underlying patterns in the data where IAGT provides better trend in design with linear and polynomial regressions.

5.6.3 Polynomial Regression with Dense Layers:

The **Polynomial Regression with Dense Layers**, in contrast, leverages regularization and gradient clipping, which helps the model avoid overfitting even when trained on a relatively small dataset. The model generates polynomial features to capture the non-linearity in the data and uses dense layers to model the interactions between the polynomial features. The gradient descent algorithm used to train the model ensures that the weights of the model are updated carefully, and **L2 regularization** penalizes large coefficients, preventing them from becoming too complex.

As observed in the evaluation, the **Polynomial Regression with Dense Layers** achieved an **R² score of 0.889** and a **Mean Squared Error (MSE) of 421.76**. While this score is slightly lower than that of **GBM (0.92)** and **Random Forest (0.912)**, the proposed model maintains a much lower risk of overfitting, making it a better generalizer. It does not require the extensive hyperparameter tuning that Decision Trees, Random Forests, and GBMs demand, especially in managing overfitting. Additionally, it doesn't suffer from the instability or complexity of ensemble methods like Random Forests or GBMs, which require careful monitoring of parameters to avoid overfitting.

5.6.4 Comparing Performance:

- The **R² score** of 0.889 for the proposed model, though slightly lower than that of **GBM** and **Random Forests**, demonstrates strong generalization capabilities. This implies that, despite being simpler than ensemble methods, the **Polynomial Regression with Dense Layers** model is able to predict the target variable accurately without overfitting to the noise in the data.
- The **MSE of 421.76** indicates that the model is making reasonably small errors in

Table-3 Representing he comparison of proposed and Existing models with performance metrics for 3k samples

Parameters	Model	MSE	R ² Score	Notes
Dataset(3k)	Linear Regression	347.69	0.909	Strong fit with less overfitting
Dataset(3k)	Decision Tree	685.34	0.821	alpha = 0.001, learning rate = 0.9
Dataset(3k)	Random Forest	337.99	0.912	Strong fit with more overfitting
Dataset(3k)	Gradient Boosting	318.75	0.92	Best performer with more over fit
Dataset(3k)	IAGT Model Phase-1	421.76	0.889	alpha = 0.01, learning rate = 0.5 with best fit no overfitting case
Dataset(3k)	IAGT Phase-2	484.85	0.887	Slight performance drop from Phase-1 with best fit no overfitting case
Dataset(3k)	IAGT Phase-3	493.76	0.87	Least performance among the IAGT phases with best fit no overfitting case

prediction, which is competitive with other models. In comparison, the **MSE for Random Forest (337.99)** and **GBM (318.75)**, although lower, indicate that while these models fit the data more closely, they may be overfitting, as evidenced by the higher variance in predictions when the model is exposed to new or unseen data.

In summary, the **Polynomial Regression with Dense Layers** performs comparably to the more complex models like **Random Forest** and **Gradient Boosting Machines**, with a minor sacrifice in performance (R² score of 0.889). However, it excels in **generalization**, avoiding the overfitting problem that often plagues more complex models, especially on smaller datasets. With **regularization** and **gradient clipping**, it ensures that model complexity remains manageable, making it a more **robust** and **efficient** choice in scenarios where overfitting is a significant concern. This makes it ideal for tasks where generalization across unseen data is critical.

5.7 TABULATIONS

The proposed IAGT model, in comparison with the existing models in table-3, demonstrates its potential as a strong performer in terms of minimizing overfitting while maintaining an optimal fit. Specifically, **IAGT Model Phase-1** achieved a Mean Squared Error (MSE) of 421.76 and an R^2 score of 0.889, outperforming many traditional machine learning models like Decision Tree (MSE = 685.34, $R^2 = 0.821$) and closely rivaling others such as **Random Forest** (MSE = 337.99, $R^2 = 0.912$) and **Gradient Boosting** (MSE = 318.75, $R^2 = 0.92$). The key differentiator of the IAGT model phases is its ability to achieve a good balance of fit without significant overfitting, as evident from its consistent performance across different phases. In contrast,

other models such as **Random Forest** and **Gradient Boosting** show strong fits but tend to exhibit higher levels of overfitting, which reduces their generalization capabilities. The IAGT model's **Phase-1** configuration with an alpha of 0.01 and learning rate of 0.5 delivers the best performance with minimal overfitting, and this behavior is consistent even in **Phase-2** (MSE = 484.85, $R^2 = 0.887$) and **Phase-3** (MSE = 493.76, $R^2 = 0.87$), where performance slightly drops but remains stable. Thus, the IAGT model provides a balanced approach for regression tasks, optimizing both model fit and generalization, particularly in settings where overfitting is a concern.

Table-4 Representing the comparison of proposed and Existing models with performance metrics for 5k samples

Parameter	Model	MSE	R^2 Score	Notes
Dataset(5k)	Linear Regression	240.69	0.919	Strong fit with less overfitting
Dataset(5k)	Decision Tree	285.34	0.881	alpha = 0.001, learning rate = 0.9
Dataset(5k)	Random Forest	237.99	0.929	Strong fit with more overfitting
Dataset(5k)	Gradient Boosting	248.75	0.931	Best performer with more over fit
Dataset(5k)	IAGT Model Phase-1	221.76	0.929	alpha = 0.01, learning rate = 0.5 with best fit no overfitting case
Dataset(5k)	IAGT Phase-2 (Dense Lasso)	284.85	0.897	Slight performance drop from Phase-1 with best fit no overfitting case
Dataset(5k)	IAGT Phase-3 (Ridge)	293.76	0.895	Least performance among the IAGT phases with best fit no overfitting case

When analyzing the combined performance of models for 5,000 samples in table-4, it is clear that the **IAGT (Incremental Adaptive Gradient Technique)** models show a significant improvement, particularly in Phase-1, when the dataset size increases. With **3,000 samples**, the IAGT models demonstrate a solid performance with no overfitting and the best generalization capabilities, with Phase-1 yielding a low MSE of 221.76 and an impressive R^2 score of **0.929**. However, when the dataset is increased to **5,000 samples**, the performance of the IAGT Phase-1 model improves even further, continuing to deliver a strong fit without overfitting, maintaining a high R^2

score of **0.929**. This improvement is evident when compared to other models like **Gradient Boosting** or **Decision Trees**, which can show increased overfitting as the sample size grows, as they attempt to fit more complex relationships in larger datasets. The **IAGT model** imparts such improvement with larger datasets lies in its design philosophy, which is crucial for dealing with the complexity of larger datasets without falling into overfitting. Unlike models that heavily rely on deep learning or complex decision trees, **IAGT** uses incremental learning and adaptive techniques that enable it to improve with each batch of data without overly fitting to any specific part of the dataset. As the number of samples

increases, the model's ability to generalize improves, and it becomes better at adjusting its parameters to capture the overall trend in the data, rather than memorizing specific data points. This is particularly beneficial as it ensures the model performs consistently well even as the training set grows in size. Therefore, designing a model like IAGT that balances complexity with generalization becomes essential for real-world applications where datasets are not only large but continuously expanding. The **IAGT Phase-1 model's** ability to adapt to larger datasets, maintain accuracy, and avoid overfitting makes it a standout choice for scalable and reliable predictive modelling.

In conclusion, the importance of designing models like **IAGT** lies in their capacity to handle increasing dataset sizes effectively, ensuring consistent and improved performance without losing generalization power. While other algorithms, such as **Gradient Boosting**, show improvement with larger datasets, they are more prone to overfitting as the complexity of the model increases. The ability of **IAGT** to improve its generalization as the sample size grows, without sacrificing accuracy, highlights the necessity for adaptable and efficient algorithms in data science, especially in real-world scenarios where datasets are continually growing and evolving.

5.7.1 Real time justification

The IAGT model, particularly when applied to cotton yield prediction, demonstrates a significant contribution to precision agriculture and sustainability. By leveraging polynomial feature generation combined with deep learning techniques, it effectively models the intricate relationships between environmental and agricultural factors such as rainfall, temperature, and soil quality, which are crucial for predicting crop yield. The model's adaptability to varying sample sizes (e.g., 3,000 and 5,000 samples) and its ability to minimize overfitting, especially when compared to traditional machine learning models like Decision Trees and Random Forests, ensures that it provides accurate and robust predictions without becoming too complex or sensitive to noise in the data. As a result, farmers can rely on the IAGT model to make data-driven decisions with higher confidence, leading to more precise management of agricultural resources. The optimization of water, fertilizer, and pesticide usage becomes more feasible, as the model can predict cotton yield with minimal error, ensuring that

resources are not overused or wasted, contributing directly to the reduction of environmental impact.

Furthermore, the IAGT algorithm's ability to generalize well across different datasets and prevent overfitting, even when dealing with larger sample sizes, makes it a powerful tool for sustainable farming. Unlike other models that may overfit and produce fluctuating performance with new data, the IAGT's performance remains stable, ensuring that its predictions are reliable even in the face of changing agricultural conditions. This helps mitigate the risk of crop loss or over-application of inputs, both of which can be costly and environmentally damaging. By improving prediction accuracy and minimizing resource waste, the IAGT model supports sustainable farming practices by providing a more efficient approach to cotton production. In the long term, as the model continues to be refined, it could help enhance overall agricultural productivity, improve yield forecasting, and assist farmers in making informed, eco-friendly decisions, ultimately leading to a more sustainable and productive agricultural system.

5.7.2 Problems and Open Research Issues

Despite the promising results of the **IAGT model** in cotton yield prediction, several challenges remain. Key issues include the availability and quality of agricultural data, especially in regions with limited infrastructure, which can affect model accuracy. Additionally, while the model adapts well over time, real-time responsiveness to sudden environmental changes such as extreme weather events or pest outbreaks needs improvement. Further research is also needed to enhance the model's interpretability and transparency, making it more accessible to farmers and agricultural experts. Moreover, scalability across diverse agricultural regions, handling noisy or incomplete data, and integrating additional factors such as pest detection and irrigation management are important areas for future exploration.

Additionally, real-world deployment and field testing in varied geographic locations are necessary to validate the model's practical utility. Research should also focus on developing cost-effective deployment strategies for smallholder farmers and conducting economic assessments to understand the model's financial impact. Addressing these issues will make the IAGT model more versatile, efficient, and accessible, facilitating its broader application in

precision agriculture and improving cotton yield predictions worldwide.

6. CONCLUSIONS

In the conclusion of this paper, the scientific contribution of the work should be clearly articulated, emphasizing how the proposed **Incremental Adaptive Gradient Technique (IAGT)** model advances the field of cotton yield prediction. This study significantly adds to the body of knowledge by addressing key challenges faced by previous prediction models, such as overfitting, poor generalization, and inadequate scalability. While traditional models like Random Forests and Gradient Boosting Machines (GBM) show high accuracy with R^2 values of 0.912 and 0.92, they tend to suffer from overfitting, especially when dealing with noisy or smaller datasets. In contrast, the **IAGT Phase-1** model achieved a competitive **R^2 score of 0.889** and an **MSE of 421.76**, illustrating its ability to generalize effectively without over-tuning, even with limited data. This showcases the strength of the **IAGT model** in adapting to real-world, noisy agricultural data without sacrificing performance, which is a significant improvement over models that are prone to overfitting.

Furthermore, the integration of **Lasso (L1)** and **Ridge (L2)** regularization techniques within the IAGT framework ensures the model controls complexity, preventing overfitting, which is critical in agricultural contexts where datasets may be small, incomplete, or noisy. This ability to avoid overfitting while maintaining strong generalization is an important scientific contribution, especially when considering the dynamic nature of agricultural environments. In Phase-1, the **IAGT model** achieved remarkable results on **larger datasets** (5,000 samples), obtaining an **R^2 score of 0.929** and **MSE of 221.76**, outperforming traditional models such as Decision Trees and Random Forests, which tend to struggle as the dataset size grows. Unlike ensemble methods, which tend to overfit as datasets expand, the **IAGT model** adapts incrementally, making it more robust and reliable for scalable predictive modeling.

The key takeaway from this research is that the **IAGT model** strikes a balance between **accuracy** and **generalization**, offering significant improvements over previous models. Its ability to adapt to increasing data and handle dynamic agricultural conditions makes it a strong candidate

for future predictive tasks in **precision agriculture**. The work's novel approach lies in its ability to integrate real-time adaptability, regularization, and scalability in cotton yield prediction, making it particularly suited for **long-term agricultural forecasting**. By addressing the **gaps** in existing models, such as the inability to effectively generalize in real-world settings, the **IAGT model** contributes valuable new knowledge to the field, setting a foundation for future advancements in **machine learning applications** in agriculture.

In summary, this paper's contribution lies in the development of a **scalable, adaptive, and regularized** prediction model for cotton yield, which demonstrates superior performance and stability across different dataset sizes and real-world conditions. This advancement not only improves cotton yield prediction but also paves the way for more reliable and adaptive predictive models in other agricultural domains, ultimately contributing to the broader goal of **enhancing food security** and **resource management** in the face of global agricultural challenges.

6.1 Scope

As dataset sizes continue to grow, the scope of applying advanced techniques like **incremental learning** and **adaptive regularization** becomes even more critical. With datasets expanding to 10,000, 12,000, 20,000, or 40,000 samples, models must evolve to handle increased computational complexity while maintaining generalization. Leveraging **mini-batch gradient descent**, **online learning**, and **distributed computing** allows for efficient training on large datasets, ensuring scalable performance. Regularization techniques, such as **L2** and **ElasticNet**, help manage model complexity and prevent overfitting, allowing for accurate predictions even as the data volume increases. This adaptability ensures robust performance and scalability in real-world predictive modeling applications.

REFERENCES:

- [1] H. Liang, R. Xie, W. Zhou, B. Li, D. Bi, and H. Chang, "Extended Grey Model Based on Particle Swarm Optimization and Its Application in Cotton Yield Prediction," 2022 4th Int. Conf. Data Intelligence Security (ICDIS), Shenzhen, China, 2022, pp. 469-473,
- [2] M. F. Celik, M. S. Isik, E. Erten, and G. Taskin, "Informative Earth Observation Variables for Cotton Yield Prediction Using Explainable

- Boosting Machine," 2023 IEEE Int. Geoscience Remote Sensing Symp. (IGARSS), Pasadena, CA, USA, 2023, pp. 3542-3545,
- [3] M. F. Celik, M. S. Isik, E. Erten, and G. Camps-Valls, "Explainability of End and Mid-Season Cotton Yield Predictors in Conus," 2023 IEEE Int. Geoscience Remote Sensing Symp. (IGARSS), Pasadena, CA, USA, 2023, pp. 3538-3541,
- [4] M. S. Isik, M. F. Celik, and E. Erten, "Interpretable Cotton Yield Prediction Model Using Earth Observation Time Series," 2023 IEEE Int. Geoscience Remote Sensing Symp. (IGARSS), Pasadena, CA, USA, 2023, pp. 3442-3445,
- [5] M. S. Isik, M. F. Celik, and E. Erten, "Unveiling the High-Resolution Cotton Yield Variations from Low-Resolution Statistics: Lessons from a Nationwide Study in Turkey," 2024 IEEE Int. Geoscience Remote Sensing Symp. (IGARSS), Athens, Greece, 2024, pp. 5040-5043,
- [6] S. T. Haider et al., "An Ensemble Machine Learning Framework for Cotton Crop Yield Prediction Using Weather Parameters: A Case Study of Pakistan," IEEE Access, vol. 12, pp. 124045-124061, 2024,
- [7] A. Mitra et al., "Cotton Yield Prediction: A Machine Learning Approach with Field and Synthetic Data," IEEE Access, vol. 12, pp. 101273-101288, 2024,
- [8] M. F. Celik, M. S. Isik, G. Taskin, E. Erten, and G. Camps-Valls, "Explainable Artificial Intelligence for Cotton Yield Prediction with Multisource Data," IEEE Geoscience Remote Sensing Lett., vol. 20, pp. 1-5, 2023, Art no. 8500905,
- [9] A. K. Uttam, "Cotton Leaves Diseases Classification Using VGG16 Based Transfer Learning," 2023 3rd Int. Conf. Ubiquitous Comput. Intell. Syst. (ICUIS), Gobichettipalayam, India, 2023, pp. 296-300,
- [10] F. Devoto et al., "Insights in the Ability of High-Resolution Narrow Band Multispectral and Thermal Sensors to Estimate Cotton Production in Australia," 2024 IEEE Int. Geoscience Remote Sensing Symp. (IGARSS), Athens, Greece, 2024, pp. 1510-1513,
- [11] A. Balmumcu, K. Kayabol, and E. Erten, "Machine Learning-based Crop Yield Prediction by Data Augmentation," 2024 32nd Signal Process. Commun. Appl. Conf. (SIU), Mersin, Turkey, 2024, pp. 1-4,
- [12] N. Agarwal, S. Ray, and K. C. Tripathi, "Time Series Forecasting of Agriculture Yield of Cotton with Regression Model Implementation," 2022 OPJU Int. Technol. Conf. Emerging Technol. Sustainable Dev. (OTCON), Raigarh, Chhattisgarh, India, 2023, pp. 1-6,
- [13] A. Srivastava, B. Singh Rawat, G. Kumar, V. Bhatnagar, and N. Garg, "Cotton Leaf Disease Prediction Using VGG16 and RESNET50," 2024 Parul Int. Conf. Eng. Technol. (PICET), Vadodara, India, 2024, pp. 1-6,
- [14] C. A. Arun, G. Giridhar, V. V. Krishna, G. V. Krishna, S. Thenappan, and M. N. Gadde, "Optimizing Plant Health Monitoring: A CNN Model Based Rice and Cotton Disease Prediction," 2023 Int. Conf. Data Sci. Agents Artif. Intell. (ICDSAAI), Chennai, India, 2023, pp. 1-6,
- [15] U. Dewangan, R. H. Talwekar, and S. Bera, "A Systematic Review on Cotton Plant Disease Detection & Classification Using Machine & Deep Learning Approach," 2023 1st DMIHER Int. Conf. Artif. Intell. Educ. Industry 4.0 (IDICAIEI), Wardha, India, 2023, pp. 1-6,
- [16] S. Mohmmad et al., "Detection and Classification of Various Diseases in Cotton Crops Using Advanced Neural Network Approaches," 2024 4th Int. Conf. Adv. Electr. Comput., Commun. Sustain. Technol. (ICAECT), Bhilai, India, 2024, pp. 1-6,
- [17] P. S. V, S. N. Raj, S. Naveen, and A. U. M, "Crop Yield Prediction of Cotton Using Optimization Technique," 2023 Int. Conf. Smart Syst. Appl. Electr. Sci. (ICSSES), Tumakuru, India, 2023, pp. 1-6,
- [18] S. Nagarajan, T. Kumaravel, P. Natesan, K. S. Nagul, A. M. Naveen, and A. Sakthisundaram, "Cotton Boll Detection Through Deep Learning Techniques," 2023 2nd Int. Conf. Autom., Comput. Renewable Syst. (ICACRS), Pudukkottai, India, 2023, pp. 1357-1362,
- [19] M. F. Shahid, S. S. Hussain, A. Zehrah, H. S. Khan, and M. A. Ahmed, "Predicting Temperature and Humidity for Cotton Field Using Deep Learning Models in Smart Agriculture System," 2023 IEEE 8th Int. Conf. Eng. Technol. Appl. Sci. (ICETAS), Bahrain, Bahrain, 2023, pp. 1-6,
- [20] A. Reyana, S. Kautish, P. M. S. Karthik, I. A. Al-Baltah, M. B. Jasser, and A. W. Mohamed, "Accelerating Crop Yield: Multisensor Data Fusion and Machine Learning for Agriculture Text Classification," IEEE Access, vol. 11, pp. 20795-20805, 2023,

- [21] S. K. Shuversa, A. A. Ryan, S. Mamun, N. Nabi, and M. S. Ahamed, "An Approach Using Machine Learning to Determine Bangladeshi Jute Yield Relying on Weather Patterns," 2022 32nd Int. Conf. Comput. Theory Appl. (ICCTA), Alexandria, Egypt, 2022, pp. 221-227,
- [22] Y. A. Madany, M. S. V, and S. S. Raj, "Cotton Disease Detection and Yield Prediction Using Machine Learning," 2023 Int. Conf. Autom. Smart Technol. (ICAST), Chennai, India, 2023, pp. 1-6,