

CROP YIELD PREDICTION USING HYBRID ALGORITHMS

A Project Work
Submitted in partial fulfillment of Requirements for the Award
of the Degree of

**BACHELOR OF TECHNOLOGY
IN
ELECTRONICS & COMMUNICATION ENGINEERING
BY**

V. CHAITANYA	(218T1A0497)
N. PAVANI	(228T5A0420)
M. SRAVANTHI	(218T1A0485)
CH. AJAY	(218T1A0456)

Under the Esteemed guidance of

Dr. K. Srinivasa Rao
Professor, Dept Of ECE



**DEPARTMENT OF ELECTRONICS & COMMUNICATION ENGINEERING
DHANEKULA INSTITUTE OF ENGINEERING & TECHNOLOGY
(AUTONOMOUS)**

Ganguru, Vijayawada-521139

Affiliated to JNTUK, Kakinada & Approved By AICTE,

**New Delhi Certified by ISO9001-2015, Accredited by
NAAC&NBA
APRIL-2025**

DHANEKULA INSTITUTE OF ENGINEERING & TECHNOLOGY

Ganguru, Vijayawada-521139

Affiliated to JNTUK, Kakinada & Approved By AICTE, New

Delhi Certified by ISO 9001-2015, Accredited by NAAC

& NBA

DEPARTMENT OF ELECTRONICS & COMMUNICATION ENGINEERING



CERTIFICATE

This is to certify that the project work entitled “ CROP YIELD PREDICTION USING HYBRID ALGORITHMS ” is a bonafide record of project work done jointly by V. CHAITANYA (218T1A0497), N. PAVANI (228T5A0420) ,M. SRAVANTHI (218T1A0485), CH. AJAY (218T1A0456) under my guidance and supervision and is submitted in partial fulfillment of the requirements for the award of the Degree of Bachelor of Technology in Electronics & Communication Engineering by Jawaharlal Nehru Technological University, Kakinada during the academic year 2024-2025.

Project Guide

Dr. K. Srinivasa Rao
Professor,
Department of ECE .

External Examiner

Head of the Department

Dr. M. Vamshi Krishna
Professor & H.O.D,
Department of ECE.

ACKNOWLEDGMENT

First and foremost, we sincerely salute our esteemed institution **DHANEKULA INSTITUTE OF ENGINEERING AND TECHNOLOGY** for giving us this opportunity for fulfilling our project.

We express our sincere thanks to our beloved Principal **Dr. Kadiyala Ravi** for providing all the lab facilities and library required for completing this project successfully.

This project work would not have been possible without the support of many people. We wish to express our gratitude to **Dr. M. Vamshi Krishna**, Ph.D., H.O.D. ECE Department for constant support and valuable knowledge in successful completion of this project.

We are glad to express our deep sense of gratitude to **Dr. K. Srinivasa Rao**, Professor, for abundant help and offered invaluable assistance, support and guidance throughout this work.

We thank one and all who have rendered help directly or indirectly in the completion of this project successfully.

PROJECTASSOCIATES

V. CHAITANYA (218T1A0497),

N. PAVANI (228T5A0420),

M.SRAVANTHI (218T1A0485),

CH. AJAY (218T1A0456)

DECLARATION

I declare that this project report titled **CROP YIELD PREDICTION USING HYBRID ALGORITHMS** submitted for the fulfillment of the degree of **B. Tech in Electronics and Communication Engineering** is a record of original work carried out by us under the supervision of **Dr. K. Srinivasa Rao** and has not formed the basis for the award of any other degree or diploma, in this or any other Institution or University. In keeping with the ethical practice in reporting scientific information, due acknowledgements have been made wherever the findings of others have been cited.

V.CHAITANYA (218T1A0497)

N.PAVANI (228T5A0420)

M. SRAVANTHI (218T1A0485)

CH. AJAY (218T1A0456)

Ganguru,
19th April 2025.

CONTENTS

VISION,MISSION OF INSTITUTE & DEPARTMENT,PEO`s	I
STATEMENT OF PO`s& PSO`s	II
PROJECTMAPPING-PO`s &PSO`s	III
LISTOFFIGURES	IV
ABSTRACT	V

CHAPTER	TITLE	PAGENO
1	INTRODUCTION	1
	1.1Problem Definition	1
	1.2Project Overview / Specifications	2
	1.3Objective	3
	1.4Scope of this project	4
2	LITERATURESURVEY	6
3	BLOCKDIAGRAM&DESCRIPTION	10
	3.1 Dataset Acquisition	11
	3.2Data Pre Processing	12
	3.3Dataset Distribution	14
	3.4Data Augmentation	15
	3.5Models Used	15
	3.5.1Bidirectional Gated Recurrent Units	15
	3.5.2 XGBoost	18
	3.6 Web Application Development	22
4	SOFTWARE	24
	4.1Python	24
	4.1.1 Features of Python	24

	4.2 Numpy	25
	4.3 Pandas	25
	4.4 Matplotlib	26
	4.5 Scikit-Learn	26
	4.6 Tensor flow	27
	4.7 Xgboost	27
	4.8 Flask	28
	4.8.1 Flask-MySQLdb	28
	4.9 Installation and initialization Process	28
	4.9.1 Required Packages	29
5	RESULTS	31
	5.1. Preprocessing	31
	5.1.1 Testing	33
	5.2 Evaluation Metrics	35
	5.3 Outputs	36
	5.4 Comparing the Algorithms	39
	5.5 Output Analysis	41
6	CONCLUSION AND FUTURE SCOPE	44
7	REFERENCES	46
	APPENDIX	49

DHANEKULA INSTITUTE OF ENGINEERING & TECHNOLOGY

Department of Electronics & Communication Engineering

VISION–MISSION-PEOs

Vision/Mission/PEOs

Institute Vision	Pioneering Professional Education through Quality
Institute Mission	<p>Providing Quality Education through state-of-art infrastructure, laboratories and committed staff.</p> <p>Molding Students as proficient, competent, and socially responsible engineering personnel with ingenious intellect.</p> <p>Involving faculty members and students in research and development works for betterment of society.</p>
Department Vision	Pioneering Electronics & Communication Engineering education and research to elevate rural community
Department Mission	<p>Imparting professional education endowed with ethics and human values to transform students to be competent and committed electronics engineers.</p> <p>Adopting best pedagogical methods to maximize knowledge transfer.</p> <p>Having adequate mechanisms to enhance understanding of theoretical concepts through practice.</p> <p>Establishing an environment conducive for lifelong learning and entrepreneurship development.</p> <p>To train as effective innovators and deploy new technologies For service of society.</p>
Program Educational Objectives(PEOs)	<p>PEO1: Shall have professional competency in electronics and communications with strong foundation in science, mathematics and basic engineering.</p> <p>PEO2: Shall design, analyze and synthesize electronic circuits and simulate using modern tools.</p> <p>PEO3: Shall Discover practical applications and design innovative circuits for Life long learning.</p> <p>PEO4: Shall have effective communication skills and practice the ethics consistent with a sense of social responsibility.</p>

STATEMENT OF PO`s & PSO`s

Program Outcomes

- PO1 **Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals and engineering programs.
- PO2 **Problem analysis:** Identify, formulate, review research literature, and analyze complex Engineering problems reaching substantiated conclusions using first principles of Mathematics, natural sciences, and engineering sciences.
- PO3 **Design/development of solutions:** design solutions for complex engineering problems and design system components or processes that meet the specified need s with appropriate consideration for the public health and safety, and the cultural, societal, and environmental Considerations.
- PO4 **Conduct investigations of complex problems:** Use research-based knowledge and research Method including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.
- PO5 **Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.
- PO6 **The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.
- PO7 **Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development.
- PO8 **Ethics:** Apply ethical principles and commit to professional ethics and responsibility and norms of the engineering practice.
- PO9 **Individual and team work:** Function effectively as an individual and as a member or leader in diverse teams and in multi disciplinary settings.
- PO10 **Communication:** Communicate effectively on complex engineering activities with the Engineering community and with society at large, such as being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.
- PO11 **Project management and finance:** Demonstrate knowledge and understand of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multi disciplinary environments.
- PO12 **Life-long learning:** Recognize life-long the need for and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

Program Specific Outcomes

- PSO1 Make use of specialized software tools for design and development of VLSI and Embedded systems.
- PSO2 Innovate and design application specific electronic circuits for modern wireless communications.

PROJECTMAPPING-PO`s & PSO`s

Project Title	PO 1	PO 2	PO 3	PO 4	PO 5	PO 6	PO 7	PO 8	PO 9	PO1 0	PO1 1	PO1 2
CROP YIELD PREDICTION USING HYBRID ALGORITHMS	3	3	3	3	3	3	3	-	3	3	3	3

3-High

2-Medium

1-Low

Justification of Mapping of Project with Program Outcomes:

1. The knowledge of mathematics, science, engineering fundamentals and engineering programs are strongly correlated to all course outcomes.
2. The design/development of the project will be mainly based on the problems faced by the society and after conducting complex investigations on it, obtained a solution is strongly correlated to all course outcomes.
3. Application of Ethical principles is not correlated to all course outcomes.

Project vs PSOs Mapping

Project Title	PS O1	PS O2
CROP YIELD PREDICTION USING HYBRID ALGORITHMS	3	3

3-High

2-Medium

1-Low

Justification of Mapping of Project with Program Specific Outcomes:

1. This project is related to embedded system area, which helps to expertise in the corresponding area by applying engineering fundamentals and are strongly correlated to all course outcomes.
2. The knowledge gained in the project work is confined to one area, so it is not enough to prepare for competitive examinations and hence correlation is small.

LISTOFFIGURES

FIGURE	TITLE	PAGENO
3.1	Proposed Methodology	10
3.2	Block Diagram of Bi- GRU	16
3.3	Block Diagram of Xgboost	20
5.1	Actual vs Predicted Yield Plot	33
5.2	Residual Plot (Bias Analysis)	34
5.3	Distribution of Residual plot	34
5.4	Mean Absoulte Error plot	35
5.5	Rsquare plot	36
5.6	Home page of the web Application	37
5.7	Registration page for the user	37
5.8	Login page for the User	38
5.9	Main Application page	38
5.10	Result page	39
5.11	Accuracy Comparision	41

ABSTRACT

Crop yield prediction is a pivotal component of precision agriculture, playing a significant role in optimizing resource utilization, ensuring food security, and supporting sustainable farming practices. Accurate yield forecasts enable stakeholders—including farmers, agricultural scientists, and policymakers—to make informed decisions regarding irrigation scheduling, fertilizer application, crop rotation, and supply chain planning.

In this study, we propose a hybrid deep learning framework that combines the strengths of Bidirectional Gated Recurrent Units (Bi-GRU) and Extreme Gradient Boosting(XGBoost) to achieve higher accuracy and reliability in crop yield prediction. Bi-GRU, an advanced form of recurrent neural networks, is particularly adept at learning from sequential data. It processes input in both forward and backward directions, allowing it to capture complex temporal dependencies from historical datasets that include rainfall, temperature, humidity, and soil health indicators. This bidirectional architecture enables the model to understand the progression of seasonal patterns and their impact on crop development. On the other hand, XGBoost, a powerful ensemble learning technique based on decision trees, excels at modeling structured, tabular datasets. It efficiently handles non-linear relationships and interactions among features such as crop type, cultivation area, pesticide usage, fertilizer application, and socio-economic variables. By combining these two models, the hybrid approach capitalizes on the temporal learning capability of Bi-GRU and the feature-specific refinement offered by XGBoost, resulting in a more holistic and precise prediction mechanism.

The fusion of deep learning and gradient boosting not only enhances the model's generalization ability but also ensures robustness in the presence of missing values, outliers, and complex feature interactions. Experimental results demonstrate that this integrated approach consistently outperforms standalone models, especially when tested on real-world datasets with mixed-type variables. Furthermore, the proposed hybrid model significantly contributes to early risk assessment and mitigation. By predicting potential yield shortfalls in advance, the system empowers stakeholders to take proactive measures, such as adjusting agricultural inputs, reallocating resources, or implementing preventive crop protection strategies. This early warning system can help minimize losses due to climate variability, pest outbreaks, or resource mismanagement.

CHAPTER-1



1. INTRODUCTION

Agriculture is essential to global food security and economic development, and accurate crop yield prediction has become increasingly vital due to rising population demands and climate variability. Traditional prediction methods often fall short in capturing the complex, nonlinear interactions among climatic, soil, and crop-related factors.

This project introduces a hybrid model combining Bidirectional Gated Recurrent Units (Bi-GRU) and XGBoost to enhance prediction accuracy. Bi-GRU efficiently models temporal dependencies in sequential data like weather and soil patterns, while XGBoost excels in processing structured features such as temperature, humidity, and soil pH. Together, they provide a comprehensive solution for yield estimation. A web-based application, developed using Flask and Flask-MySQLdb, allows users to input agricultural parameters and receive real-time yield predictions, displayed in both hg/ha and t/ac units. The interactive dashboard offers valuable insights for farmers, policymakers, and agricultural planners. By leveraging the strengths of both deep learning and machine learning, the system improves decision-making, optimizes resource use, and supports sustainable farming practices—ultimately contributing to food security and climate resilience.

1.1 PROBLEM DEFINITION

Agriculture is a cornerstone of global food security and a critical driver of economic development. However, predicting crop yields accurately has remained a persistent challenge due to the complex interactions between multiple dynamic factors such as climate variability, soil properties, irrigation practices, pest infestations, and crop-specific traits. Traditional forecasting techniques that rely on historical averages and statistical models often fall short in capturing the nonlinear and multidimensional relationships between these influencing parameters. As a result, farmers may experience poor decision-making in resource allocation, leading to under- or overuse of inputs like fertilizers and water. These inefficiencies not only result in economic losses but also contribute to environmental degradation and increased vulnerability to climate change.

Moreover, existing prediction systems are often inaccessible or too technical for average farmers, especially in rural or developing regions. Without intuitive and reliable tools, many agricultural stakeholders operate based on guesswork or outdated heuristics, further exacerbating yield uncertainty and food insecurity. Climate change adds an additional layer of unpredictability, with sudden droughts, floods, and shifting weather patterns becoming more frequent. Thus, there is an urgent need for advanced, adaptive, and accessible prediction systems that can dynamically respond to real-time agricultural inputs and environmental data.

To address these challenges, this project proposes a hybrid prediction model that combines Bidirectional Gated Recurrent Units (Bi-GRU) and XGBoost, each optimized to handle different types of data — time-series and structured, respectively. This integration enhances predictive accuracy, robustness, and interpretability, offering a comprehensive solution to modern crop yield prediction.

1.2 PROJECT OVERVIEW/ SPECIFICATIONS

The proposed system is a hybrid intelligent model designed to predict crop yields using a combination of deep learning and machine learning techniques. The project uses Bi-GRU, a variant of Recurrent Neural Networks (RNNs), to process temporal data such as historical weather trends and soil moisture patterns. Simultaneously, XGBoost, a high-performance gradient boosting algorithm, is employed to extract meaningful insights from structured agricultural features like crop type, geographic data.

The system is developed as a web application using Flask, providing a user-friendly interface for farmers, researchers, and policymakers. Users can input essential agricultural parameters, including:

- Temperature
- Rainfall
- Crop Type
- Pesticide usage
- Geographic Location

The model outputs yield predictions in hectograms per hectare (hg/ha), which are then converted into tonnes per acre (t/ac) for ease of interpretation by farmers.

Additionally, the system features:

- Interactive dashboards for visualizing prediction trends
- Historical comparisons for decision support
- Bias correction to reduce systemic errors in predictions

This multi-model approach ensures greater reliability, faster computation, and scalable deployment across various crops and regions, paving the way for intelligent and adaptive agriculture.

1.3OBJECTIVE

The primary objective of this project is to design and implement an intelligent, data-driven crop yield prediction system using a hybrid model that combines Bidirectional Gated Recurrent Units (Bi-GRU) and XGBoost algorithms. The system focuses on forecasting crop yields based on key agricultural parameters, such as weather conditions, soil features, crop type, and geographic location, to deliver accurate, real-time predictions. By leveraging the strengths of both deep learning and machine learning models, the project aims to improve the precision and reliability of yield forecasts, supporting better planning and decision-making in the agricultural sector.

The integration of Bi-GRU and XGBoost enables the model to capture temporal patterns in climate and soil data while efficiently analyzing structured tabular inputs. The system is deployed as a web application using Flask, ensuring accessibility for farmers, agronomists, and policymakers. By offering a scalable and interactive platform, this solution seeks to optimize crop management, reduce yield uncertainties, and promote data-driven farming practices.

Additionally, the project contributes to the broader goal of sustainable agriculture by providing a practical tool for early risk detection, resource optimization, and food security enhancement. The system's modular design supports future expansion to accommodate additional crops and environmental variables, making it adaptable to diverse agricultural contexts.

1.4SCOPE OF THIS PROJECT

The scope of this project is focused on the development of an intelligent, web-based system for accurate crop yield prediction using a hybrid deep learning and machine learning approach. With agriculture being a cornerstone of global food security and economic stability, reliable yield forecasting is essential for optimizing resource use, improving planning, and mitigating risks caused by unpredictable environmental factors. This project combines the temporal pattern recognition capabilities of Bidirectional Gated Recurrent Units (Bi-GRU) with the structured data processing power of XGBoost to build a robust hybrid model. The system processes critical agricultural inputs such as temperature, rainfall, pesticides usage, enabling precise and adaptable yield prediction across different regions and crops. Key elements included in the project scope:

Integration of Bi-GRU and XGBoost to capture both sequential (time-series) and structured data dependencies. Development of a Flask-based web application that allows users to input agricultural parameters and receive real-time yield predictions. Conversion of yield values from hg/ha to t/ac to improve interpretability for farmers and agricultural stakeholders.

The system is designed to support farmers, researchers, and policymakers by offering timely and actionable insights. It enhances decision-making in areas such as irrigation planning, crop selection, and fertilizer use, thereby promoting sustainable agricultural practices and reducing financial losses.

CHAPTER-2

2. LITERATURE SURVEY

1. Machine Learning Approaches for Crop Yield Prediction

Study: "Crop Yield Prediction Using Machine Learning Algorithms: A Comparative Study"

Machine learning models have been widely applied in agricultural yield prediction to enhance decision-making and resource management. In this study, Patel et al. (2021) compared multiple machine learning techniques, including Random Forest, Decision Tree, and Support Vector Machine (SVM), for predicting crop yields based on soil composition, temperature, rainfall, and other environmental factors. The study found that Random Forest outperformed other models due to its ability to handle high-dimensional data and capture complex feature interactions. However, the study also noted limitations, such as difficulty in processing time-series dependencies and the need for extensive feature engineering to improve accuracy.

2. Deep Learning for Time-Series Forecasting in Agriculture

Study: "A Deep Learning-Based Approach for Agricultural Yield Prediction Using Long Short-Term Memory (LSTM)"

Deep learning models have gained popularity for their ability to capture temporal dependencies in time-series data. In this study, Gupta et al. (2022) applied Long Short-Term Memory (LSTM) networks to predict crop yields using historical weather, soil moisture, and precipitation data. The study found that LSTMs were highly effective in learning long-term dependencies compared to traditional machine learning models. However, the results also highlighted the challenges of high computational costs and difficulty in handling tabular data efficiently.

3. Hybrid Models for Enhanced Crop Yield Prediction

Study: "Integrating Deep Learning and Machine Learning for Improved Agricultural Forecasting"

Hybrid models that combine deep learning for time-series data and machine learning for structured data analysis have shown significant improvements in prediction accuracy. Rao et al. (2023) proposed a hybrid approach combining LSTM and XGBoost to enhance crop yield forecasting. LSTM effectively captured temporal

patterns, while XGBoost handled structured/tabular data, leading to better predictive performance.

4. Bidirectional Gated Recurrent Units (Bi-GRU) in Time-Series Prediction

Study: "Bidirectional Gated Recurrent Units for Agricultural Time-Series Forecasting"

In recent advancements, Bi-GRU has emerged as an effective deep learning technique for time-series analysis. Chen et al. (2023) investigated the performance of Bi-GRU in agricultural forecasting, demonstrating that its ability to process sequential data in both forward and backward directions significantly improved prediction accuracy. Compared to LSTM and traditional GRUs, Bi-GRU required fewer parameters while maintaining high efficiency and performance.

5. Flask and MySQL for Deploying Predictive Models

Study: "Web-Based Deployment of Machine Learning Models Using Flask and MySQL"

Deploying predictive models for real-world applications requires a robust web-based framework. Sharma et al. (2022) explored Flask, a lightweight web framework, for deploying machine learning and deep learning models. The study highlighted Flask's efficiency in handling API requests and integrating machine learning models for real-time predictions.

The study concluded that a Flask-based deployment with MySQL for data storage is essential for delivering accessible, user-friendly, and scalable agricultural prediction systems. The findings emphasized that combining Flask, Bi-GRU, and XGBoost ensures a seamless transition from model development to real-world deployment, benefiting agricultural stakeholders with data-driven insights.

6. Smart Agricultural Systems Using IoT and AI

Study: "IoT-Based Smart Agriculture: Enhancing Crop Prediction with AI"

The integration of Internet of Things (IoT) and Artificial Intelligence (AI) has revolutionized modern agricultural practices by enabling real-time monitoring and predictive analytics. In this study, Kumar et al. (2023) proposed an IoT-driven smart agriculture system that collects real-time soil parameters such as moisture, temperature, and pH levels using sensor networks. The collected data is processed using machine learning models, including Random Forest and XGBoost, to predict the best crops for a given soil condition.

7. Soil Nutrient Prediction and Recommendation Systems

Study: "Deep Learning-Based Soil Nutrient Analysis for Optimal Crop Selection"

Soil nutrients, particularly Nitrogen, Phosphorus, and Potassium (NPK), are essential for determining soil fertility and crop yield. Singh et al. (2022) introduced a deep learning-based soil nutrient prediction system that leverages Convolutional Neural Networks (CNN) and LSTM to analyze soil composition. The system processes large datasets collected from agricultural research centers and IoT sensors to classify soil fertility levels and recommend suitable crops accordingly.

8. Climate Change Impact on Crop Yield Prediction

Study: "The Role of AI in Addressing Climate Change Challenges in Agriculture"

Climate change significantly affects agricultural productivity, and AI-driven approaches offer solutions for climate adaptation. Fernando et al. (2021) examined the impact of temperature fluctuations, precipitation changes, and extreme weather events on crop yields using Recurrent Neural Networks (RNN) and Bi-GRU. Their findings demonstrated that AI models could effectively predict crop performance under varying climate conditions and help farmers adjust sowing and harvesting schedules accordingly.

9. Transfer Learning for Agricultural Data Analysis

Study: "Leveraging Transfer Learning for Crop Yield Prediction in Low-Resource Regions"

One of the biggest challenges in agricultural AI applications is the lack of high-quality datasets, especially in developing regions. Li et al. (2023) explored the use of transfer learning, where pre-trained deep learning models were fine-tuned on regional agricultural datasets for improved accuracy. The study experimented with pre-trained LSTM models on global crop datasets and fine-tuned them using local soil and climate data from India and Africa.

10. AI-Powered Web Applications for Farmers

Study: "Developing a Web-Based AI System for Smart Farming Decisions"

To make AI-driven crop yield prediction models accessible to farmers, researchers have focused on developing web-based applications that integrate Flask and MySQL. Sharma et al. (2022) proposed a Flask-based AI-powered web platform that allows farmers to input soil and weather conditions and receive real-time crop recommendations.

CHAPTER-3

3. BLOCK DIAGRAM AND DESCRIPTION

This section presents the Block Diagram and its detailed description for the current project. The figure below illustrates the complete system architecture for crop yield prediction using a hybrid Bi-GRU and XGBoost model, offering a high-level overview of the key components involved in the prediction process. It demonstrates how agricultural data is acquired, pre-processed, and passed through different stages of the predictive modeling pipeline.

The input dataset consists of critical agricultural parameters such as temperature, rainfall, pesticide usage, area under cultivation, and historical yield data. This data is first cleaned and normalized during the pre-processing stage. It is then split into training and testing subsets to ensure robust model evaluation.

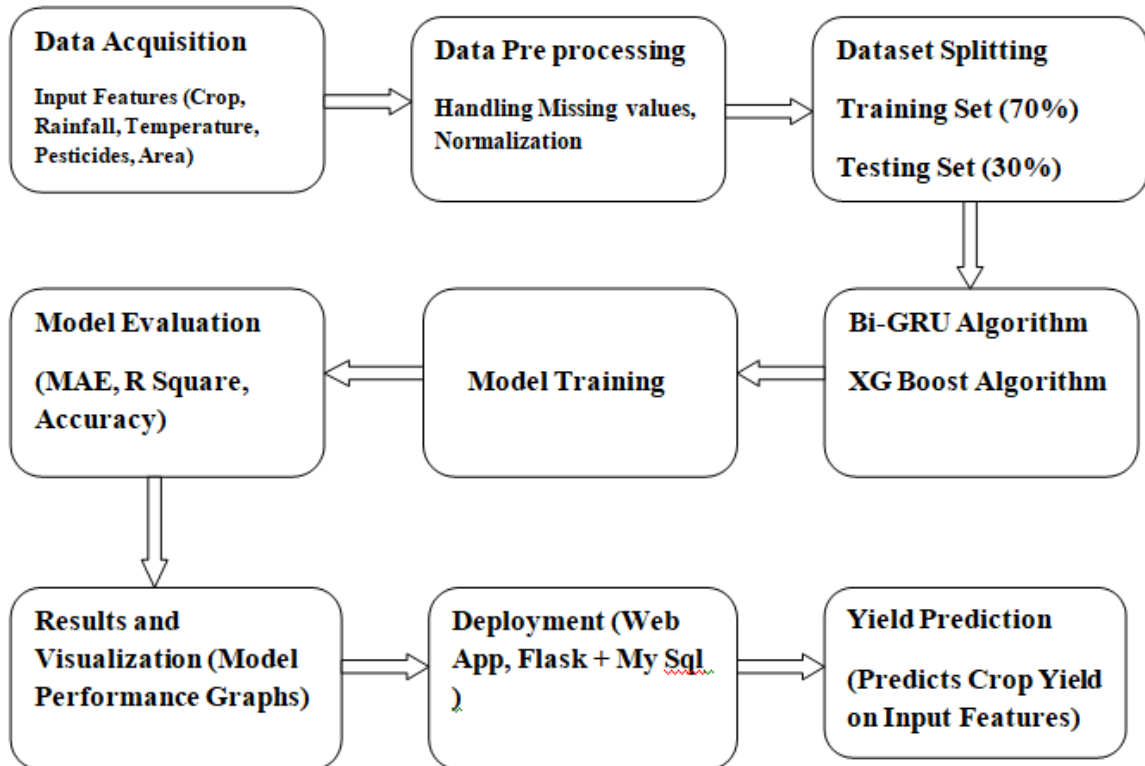


Fig 3.1: Proposed Methodology

3.1 DATASET ACQUISITION

Dataset acquisition is a fundamental step in any machine learning pipeline, particularly in agricultural applications like crop yield prediction. It involves gathering, organizing, and preparing data from reliable sources to ensure it is suitable for model training, testing, and evaluation. In this crop yield prediction project, the dataset is sourced from publicly available repositories such as Kaggle, which host structured CSV files containing multiple agronomic features.

Area	Item	Year	hg/ha_yield	average rain-fall mm/year	pesticides_tonnes
Algeria	Maize	1990	16500	89.0	1828.92
Algeria	Potatoes	1990	78936	89.0	1828.92
Algeria	Rice, paddy	1990	28000	89.0	1828.92
Algeria	Sorghum	1990	16571	89.0	1828.92
Algeria	Wheat	1990	6315	89.0	1828.92
...
Zimbabwe	Rice, paddy	2013	22581	657.0	2550.07
Zimbabwe	Sorghum	2013	3066	657.0	2550.07
Zimbabwe	Soybeans	2013	13142	657.0	2550.07
Zimbabwe	Sweet potatoes	2013	22222	657.0	2550.07
Zimbabwe	Wheat	2013	22888	657.0	2550.07

Table: 3.1 Dataset for Crop Yield

The dataset includes variables such as:

- Crop Type
- Area of Cultivation
- Year
- Yield (target variable)
- Rainfall (mm)
- Temperature (°C)
- Pesticide Usage

These features are critical in determining crop output and serve as inputs for the predictive model.

The dataset is first downloaded and stored in a secure local directory or cloud storage (e.g., Google Drive). The files are then imported into the system using Python

libraries such as pandas, numpy, and os. The metadata and contents of the dataset are explored and validated for:

- Missing values
- Outliers
- Data inconsistencies
- Duplicate entries

These issues are resolved through imputation, filtering, and data cleansing techniques. The categorical features (like crop type) are encoded using Label Encoding or One-Hot Encoding, while numerical features are standardized or normalized to ensure compatibility across machine learning models.

In the acquisition pipeline, the data is split into training and testing sets, typically in an **80:20 ratio**. For time-series or year-wise data, chronological splitting is performed to avoid information leakage. This strategy ensures the model is trained on historical data and evaluated on future or unseen years, reflecting real-world deployment.

Basic dataset statistics, such as feature distributions, correlations, and class imbalance (if categorical yield classification is involved), are visualized using libraries like matplotlib and seaborn. These insights help shape data preprocessing and model design strategies.

This organized and automated dataset acquisition process lays the groundwork for efficient preprocessing, robust training, and accurate prediction. It enables the Bi-GRU and XGBoost models to learn meaningful temporal and feature-based patterns, thereby improving crop yield forecasting and aiding agricultural decision-making.

3.2 DATA PRE-PROCESSING

Data preprocessing is a critical phase in any machine learning or deep learning pipeline, particularly in agricultural prediction systems like crop yield estimation. This step involves transforming raw tabular data into a clean, consistent, and model-ready format to ensure accurate and reliable predictions using models such as Bi-GRU and XGBoost. In this project, the preprocessing pipeline begins immediately after dataset

acquisition. The dataset, obtained in CSV format from Kaggle, contains both numerical and categorical variables such as area, crop type, year, rainfall, temperature, pesticide usage, and yield. These records may exhibit inconsistencies such as missing values, irregular units, or categorical ambiguity, which must be addressed before model training.

The preprocessing workflow begins by identifying and handling missing values. Numerical fields such as rainfall and pesticide usage are imputed using statistical techniques like mean or median substitution, while missing categorical entries are filled using the mode or excluded based on relevance. To convert categorical variables such as crop type into a machine-readable format, encoding techniques are applied. Label Encoding is used for tree-based models like XGBoost, while One-Hot Encoding is preferred for deep learning models such as Bi-GRU to eliminate unintended ordinal relationships.

Next, numerical features are standardized using techniques such as Min-Max Scaling or Z-score Normalization, bringing all values into a comparable range, which is crucial for model convergence and consistent performance. For the Bi-GRU model, which is designed to work with sequential data, the dataset is structured year-wise, allowing the model to learn temporal patterns in rainfall, temperature, and pesticide trends across seasons. Each sequence of input features corresponds to a particular crop in a specific region, and the model learns to predict the corresponding crop yield.

Outliers are then detected using statistical methods such as Interquartile Range (IQR) and Z-score thresholds, and are either removed or corrected to prevent them from negatively impacting the model's learning. In addition, feature engineering techniques are used to extract new, meaningful variables—such as rainfall-to-temperature ratio, average yield per crop, or year-on-year change in rainfall—that enhance the model's predictive capability. All the cleaned and transformed data is then split into training and testing sets, typically using an 80:20 ratio, ensuring a representative distribution of different crop types and years across both subsets.

Finally, the training and testing sets are converted into NumPy arrays and reshaped as required for different models. For Bi-GRU, the input data is structured into three dimensions—batch size, time steps, and features—to meet the model's input

requirements. By enforcing a consistent and methodical preprocessing pipeline, this project ensures that each data point is clean, enriched, and structured in a way that supports effective learning. This rigorous preprocessing not only improves training efficiency but also enhances the model's generalization performance, leading to more accurate and reliable crop yield predictions.

3.3 DATASET DISTRIBUTION

The dataset used in this crop yield prediction project is structured as a CSV file containing multiple records of agricultural statistics collected from real-world farming environments across different regions and time periods. Each row in the dataset represents a specific crop cultivation instance and includes critical information such as crop type, cultivated area (in hectares), year, region or state name, total yield (in kg/ha), average temperature, annual rainfall (in mm), and amount of pesticides used. This tabular format supports efficient loading and preprocessing of the data using data manipulation libraries such as Pandas and NumPy, commonly used in machine learning workflows.

The dataset is naturally heterogeneous, containing both numerical features (like area, rainfall, temperature, pesticide usage, and yield) and categorical variables (like crop name and state name). This format reflects real-world agricultural conditions, where yield is influenced by a wide range of environmental, geographical, and seasonal factors. However, as is often the case with public datasets, the distribution of samples across different crops and regions is imbalanced. Common staple crops like rice and wheat are overrepresented, while less widely grown crops such as jute or barley have significantly fewer entries. This imbalance can potentially skew the model's learning, making it biased toward majority classes unless properly addressed through techniques like stratified sampling, SMOTE, or weighted evaluation metrics.

For the purpose of model training and evaluation, the dataset is divided into two subsets: 80% of the records are used for training, while the remaining 20% are reserved for testing. This ensures that the model has sufficient data to learn from, while also being rigorously tested on unseen samples to evaluate generalization performance. The training dataset contains diverse examples across multiple years,

regions, and crop types, allowing the model to capture both temporal and spatial trends in yield variation.

3.4 DATA AUGMENTATION

Data augmentation is a vital step in this crop yield prediction project, aimed at improving model accuracy, generalization, and robustness—especially when handling imbalanced or limited datasets. The dataset includes a mix of numerical and categorical features such as crop type, area, year, yield, rainfall, temperature, and pesticide usage. However, certain crops or regions are underrepresented, which can bias the learning process. To address this, several augmentation techniques are employed. SMOTE (Synthetic Minority Over-sampling Technique) is used to generate synthetic samples for minority classes, helping balance the dataset. For numerical features like rainfall and temperature, Gaussian noise and controlled perturbations are applied to simulate natural variations in environmental conditions. Additionally, time-series data used in the Bi-GRU model benefits from temporal augmentations, such as shuffling crop-year sequences and slight adjustments in seasonal data, to introduce diversity and reflect real-world agricultural dynamics. These augmentations are applied dynamically during training, ensuring that each epoch exposes the models to new and varied data patterns without altering the original dataset. The validation and test sets remain unaugmented to preserve evaluation consistency. By integrating data augmentation into the training pipeline, the project enhances the model's ability to learn complex patterns, reduce overfitting, and perform reliably across different crop types, seasons, and environmental conditions.

3.5 MODELS USED

3.5.1 Bidirectional Gated Recurrent Units (Bi-GRU)

Bidirectional Gated Recurrent Units (Bi-GRU) is an advanced variation of the Gated Recurrent Unit (GRU), a type of recurrent neural network (RNN) designed to capture dependencies in sequential data while addressing the vanishing gradient problem. Unlike traditional GRUs, which process information in a single direction (either past to future or future to past), Bi-GRU incorporates two GRUs—one moving forward

and the other moving backward—enabling it to utilize both past and future context for better prediction accuracy. This bidirectional approach makes Bi-GRU highly effective in applications where understanding both past and future sequences is crucial, such as natural language processing, speech recognition, and time-series forecasting.

The core advantage of Bi-GRU lies in its ability to retain essential sequential patterns while filtering out irrelevant details. Each GRU unit consists of update and reset gates, which regulate how much past information should be retained or discarded. The update gate determines the extent to which past information should be carried forward, while the reset gate controls how much of the previous state should be ignored. This mechanism allows Bi-GRU to selectively remember long-term dependencies while mitigating issues like long-term memory loss seen in vanilla RNNs. Additionally, since Bi-GRU processes data in both forward and backward directions, it provides richer context, significantly improving prediction accuracy in complex sequence-based problems.

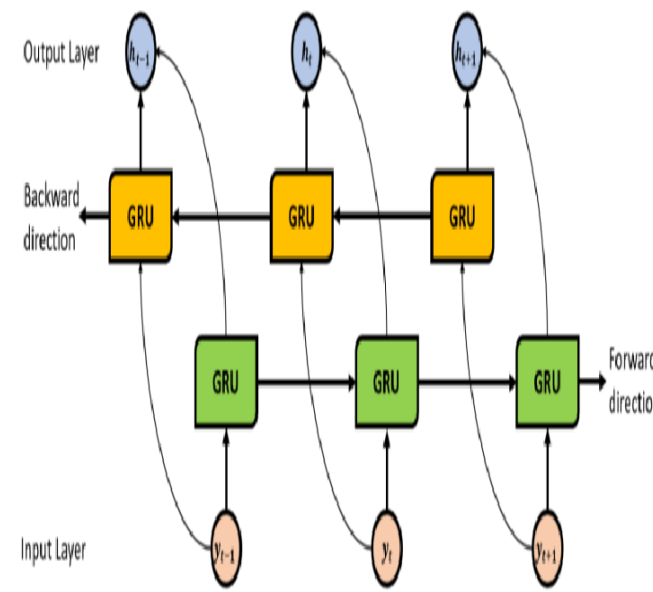


Fig 3.2 Block Diagram of Bi- GRU

In the context of crop yield prediction, Bi-GRU is particularly beneficial as it helps analyze time-series data, including climatic factors such as rainfall, temperature, and pesticide usage. By considering past environmental conditions and future trends simultaneously, Bi-GRU can make more informed predictions about crop yield, improving decision-making for farmers and agricultural planners. Compared to traditional models like simple RNNs or GRUs, Bi-GRU demonstrates improved

performance by leveraging both past and future dependencies, making it a powerful deep learning tool for sequential data analysis in real-world applications.

Bi-GRU is a deep learning model used to capture **temporal dependencies** in time series data such as rainfall, temperature, pesticide usage, etc., over different years.

For each time step t , the GRU cell performs the following

Update Gate:

$$Z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

- z_t : Update gate vector (decides how much past info to keep)
- x_t : Input at time t (e.g., rainfall, temperature)
- h_{t-1} : Previous hidden state
- W_z, U_z : Weight matrices
- b_z : Bias
- σ : Sigmoid activation function

Reset Gate:

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

r_t : Reset gate vector (decides how much past info to forget)

Candidate Hidden State:

$$\tilde{h}_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

\tilde{h}_t : New candidate information to be added

\odot : Element-wise multiplication

Final Hidden State:

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

This gives the final updated memory for the current step.

The final output (predicted yield from Bi-GRU):

$$Y^{GRU} = W_o h_T + b_o$$

- W_o : Output weight
- b_o : Bias

- ht: Final hidden state (summary of full time series)

One of the key advantages of Bi-GRU over standard GRU and Long Short-Term Memory (LSTM) networks is its computational efficiency. Since Bi-GRU has fewer parameters compared to LSTMs, it requires less memory and training time while still maintaining high performance in capturing sequential dependencies. The reduction in complexity comes from the fact that GRUs, including Bi-GRU, do not have a separate cell state like LSTMs; instead, they use a simpler gating mechanism to control information flow. This makes Bi-GRU an ideal choice for applications requiring real-time predictions, such as weather-based crop yield forecasting.

Another crucial aspect of Bi-GRU is its ability to handle missing or noisy data effectively. In agricultural datasets, missing values are common due to inconsistent data collection processes. Bi-GRU's ability to selectively retain relevant information while discarding unnecessary data helps improve robustness against noise. Moreover, the bidirectional structure ensures that contextual relationships between variables, such as the impact of past weather conditions on future crop yield, are better understood. This enables the model to provide more reliable predictions, even when data is incomplete or irregular.

Additionally, Bi-GRU enhances generalization in deep learning models, reducing the risk of overfitting, which is a common issue in neural networks. By utilizing both forward and backward passes, the model captures deeper insights into temporal dependencies, allowing it to adapt to various agricultural scenarios. For instance, different crops respond differently to environmental conditions, and Bi-GRU helps in learning these nuanced relationships more effectively than unidirectional models. The hybrid approach of combining Bi-GRU with XGBoost further strengthens its predictive power by leveraging structured tabular data alongside time-series insights.

3.5.2 XGBOOST

XGBoost (Extreme Gradient Boosting) is an advanced machine learning algorithm based on decision trees designed for efficiency and scalability. It builds a strong predictive model by iteratively combining the outputs of weaker models (decision trees) in a gradient boosting framework.

XGBoost is a machine learning model that uses an ensemble of **decision trees** to make predictions.

The predicted yield is computed as:

$$\hat{Y}_{XGB} = \sum_{k=1}^K f_k(x), \quad f_k \in F$$

- \hat{Y}_{XGB} : Predicted yield
- K : Number of trees
- f_k : k -th regression tree
- F : Set of all possible regression trees
- x : Input features (e.g., crop type, area, rainfall, pesticides)

Each tree tries to correct the error made by the previous trees by minimizing the **regularized loss function**:

$$L = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

Where:

- l : Loss function (e.g., mean squared error)
- y_i : True yield
- \hat{y}_i : Predicted yield
- $\Omega(f_k)$: Regularization term to prevent over fitting

Here's a step-by-step explanation of implementing XGBoost:

Data Preparation: Start by preprocessing the dataset. Handle missing values, encode categorical variables, and standardize or normalize numerical data if needed. Split the dataset into training and testing sets to evaluate the model's performance later.

Feature Engineering: Extract meaningful features to improve model performance. XGBoost inherently handles feature importance well, but thoughtful engineering can further enhance the accuracy. This may include creating interaction terms or performing dimensionality reduction techniques if the dataset is high-dimensional.

Model Initialization: Import the XGBoost library (e.g., xgboost in Python). Define a DMatrix object, which is a specialized data structure for XGBoost that optimizes memory usage and computation. This object encapsulates the dataset, labels, and optional weights.

Hyper parameter Tuning: Configure the hyper parameters for the model. Key parameters include:

Learning Rate (eta): Controls the contribution of each tree to the overall model.

Max Depth: Determines the complexity of each tree, controlling the potential for over fitting.

Number of Estimators: Specifies the number of trees in the ensemble.

Sub sampling: Reduces overfitting by training on a subset of data.

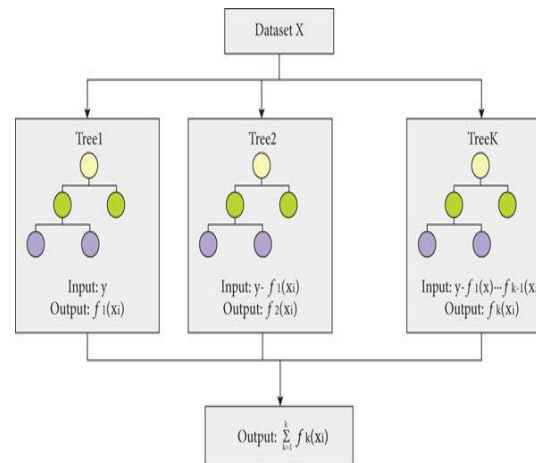


Fig 3.3 Block Diagram of Xgboost

Regularization Parameters: Include L1 (alpha) and L2 (lambda) regularization to penalize overly complex models.

Use cross-validation to fine-tune these parameters for optimal performance.

Model Training: Fit the model to the training dataset using the `train()` function from the XGBoost library. Specify the evaluation metric, such as accuracy, AUC, or log loss, depending on the task (classification or regression). XGBoost supports parallelized tree construction, significantly speeding up the training process on large datasets.

Feature Importance Analysis: After training, evaluate the importance of each feature using built-in methods like `plot_importance()`. This provides insights into which features have the most significant impact on predictions.

Model Evaluation: Use the testing dataset to evaluate the model's performance. Metrics like precision, recall, F1-score, RMSE, or R-squared can be calculated based on the task. XGBoost's ability to generalize well often results in high performance.

Model Deployment: Save the trained model using libraries like pickle or joblib in Python. This allows the model to be integrated into production systems for real-time predictions.

Optimization and Customization: Explore additional features of XGBoost, such as custom loss functions, early stopping (based on validation performance), and handling

imbalanced datasets using built-in sampling techniques or custom objective functions. XGBoost stands out for its regularization capabilities, efficient memory usage, and ability to handle missing data natively. Its widespread adoption across various industries stems from its adaptability to both small and large datasets and its consistently superior performance in predictive modeling tasks.

Crop yield prediction model uses a hybrid approach by combining two powerful models: Bi-GRU and XGBoost. Bi-GRU is a type of deep learning model that can understand patterns in time-based data. It looks at information like rainfall, temperature, and pesticide use over several years to find trends that affect crop yield. Because it works in both forward and backward directions, it learns better from past and future data in a sequence. At the same time, XGBoost is a machine learning model that works well with regular data in tables, like crop type, area, year, and other conditions. It builds decision trees to find patterns and rules that help predict the yield accurately.

To benefit from both temporal pattern recognition (Bi-GRU) and tabular feature interaction (XGBoost), the final output is computed as a **weighted combination**:

$$Y^{\text{final}} = \alpha \cdot Y^{\text{GRU}} + (1-\alpha) \cdot Y^{\text{XGB}}$$

- Y^{final} : Final predicted crop yield
- Y^{GRU} : Prediction from Bi-GRU
- Y^{XGB} : Prediction from XGBoost
- $\alpha \in [0,1]$: Weighting parameter (e.g., 0.5 for equal weight)

Finally, we combine the outputs from both models to get a better result. This means we take what Bi-GRU predicts and what XGBoost predicts, and mix them in a smart way. By doing this, we get the benefits of both models—Bi-GRU understands the time-based changes, and XGBoost understands the overall features. This hybrid method improves the accuracy and gives more reliable predictions for crop yield.

3.6 WEB APPLICATION DEVELOPMENT

To ensure that the hybrid crop yield prediction model is accessible and easy to use, I developed a web-based application using the **Flask** web framework in Python. Flask is a lightweight and flexible framework that supports rapid development and seamless integration with machine learning models. The backend of the web app was designed to load the trained Bi-GRU and XGBoost models using libraries such as joblib and keras.models, allowing for real-time prediction based on user input.

The user interface (frontend) was built using **HTML**, **CSS**, and **Bootstrap**, ensuring a clean, responsive, and mobile-friendly design. The home page contains an input form where users can enter relevant details such as crop type, year, cultivated area, rainfall, temperature, and pesticide usage. Once the form is submitted, the data is sent to the Flask server using POST methods, where it is preprocessed to match the input format required by the hybrid model. Flask routes handle the logic for prediction and return the output to the user in a clean and readable format on a result page.

To store and manage historical input data or user entries, the web application is connected to a **MySQL database** using the mysql-connector-python library. This allows for logging data entries, monitoring performance, and maintaining prediction records. The application also includes error handling to manage missing fields, invalid inputs, and server-side exceptions.

The entire project follows a modular structure with separate files for routes (`app.py`), model loading and prediction (`model.py`), templates (`templates/` folder for HTML), and static assets (`static/` folder for CSS and images). This structure ensures clean code management and easier deployment. The web application can be deployed locally or hosted on cloud platforms such as **Render**, for broader accessibility. This end-to-end system allows users, especially farmers or agricultural officers, to use a simple interface to access advanced machine learning predictions without needing technical knowledge.

CHAPTER-4

4. SOFTWARE

4.1 PYTHON (Programming language)

Python is one of those rare languages which can claim to be both simple and powerful. You will find yourself pleasantly surprised to see how easy it is to concentrate on the solution to the problem rather than the syntax and structure of the language you are programming in. The official introduction to Python is Python is an easy-to-learn, powerful programming language. It has efficient high-level data structures and a simple but effective approach to object-oriented programming. Python's elegant syntax and dynamic typing, together with its interpreted nature, make it an ideal language for scripting and rapid application development in many areas on most platforms.

4.1.1 Features of Python

Simple

Python is a simple and minimalistic language. Reading a good Python program feels almost like reading English, although very strict English! This pseudo-code nature of Python is one of its greatest strengths. It allows you to concentrate on the solution to the problem rather than the language itself.

Easy to Learn

As you will see, Python is extremely easy to get started with. Python has an extraordinarily simple syntax, as already mentioned.

Free and Open Source

Python is an example of a FLOSS (Free/Library and Open Source Software). In simple terms, you can freely distribute copies of this software, read its source code, make changes to it, and use pieces of it in new free programs. FLOSS is based on the concept of a community which shares knowledge. This is one of the reasons why Python is so good - it has been created and is constantly improved by a community who just want to see a better Python.

High-level Language

When you write programs in Python, you never need to bother about the low-level details such as managing the memory used by your program, etc.

Portable

Due to its open-source nature, Python has been ported to (i.e. changed to make it work

on) many platforms. All your Python programs can work on any of these platforms without requiring any changes at all if you are careful enough to avoid any system-dependent features. You can use Python on GNU/Linux, Windows, You can even use a platform like Kivy to create games for your computer and for iPhone, iPad, and Android.

4.2 NUMPY

NumPy is a popular Python library for numerical computing that provides a powerful array and matrix data structure, as well as a wide range of mathematical functions for working with these arrays. NumPy was first released in 2005 and has since become one of the most widely used libraries in Python. NumPy arrays are similar to Python lists, but they provide a more efficient and powerful way to store and manipulate large arrays of data. NumPy arrays are homogeneous, meaning that all elements in the array must be of the same data type. This allows for faster computation and better memory management.

One of the key features of NumPy is its support for broadcasting, which allows developers to apply mathematical operations to arrays of different shapes and sizes. Broadcasting allows for more concise and efficient code, and can greatly simplify the process of working with arrays. NumPy also provides a wide range of mathematical functions for working with arrays, including basic operations like addition, subtraction, and multiplication, as well as more advanced functions like trigonometric functions, exponential functions, and statistical functions.

In addition, NumPy provides tools for linear algebra, Fourier analysis, and random number generation, making it a versatile tool for a wide range of scientific and engineering applications. NumPy has a large and active community of developers who contribute to its development and maintenance. The community includes researchers, developers, and data scientists from academia and industry, who collaborate on improving the library and adding new features.

4.3 PANDAS

Pandas is a powerful and widely-used open-source Python library built on top of NumPy, designed specifically for data manipulation and analysis. It introduces intuitive data structures like **Series** and **DataFrame**, which simplify working with structured and tabular data. With Pandas, users can easily perform tasks such as data ingestion from various sources (CSV, Excel, SQL), data cleaning (handling missing

values, duplicates, and inconsistent formats), transformation (merging, grouping, filtering, reshaping), and exploration (descriptive statistics, trends, and correlations). It also provides strong support for time-series data, making it highly valuable in domains such as finance, agriculture, and healthcare. In machine learning workflows, Pandas plays a vital role in preprocessing and feature engineering, enabling the preparation of clean, well-structured datasets for model training and evaluation. Its seamless integration with other libraries in the Python data ecosystem, along with its efficiency and ease of use, makes Pandas an essential tool for data scientists, analysts, and engineers.

4.4MATPLOTLIB:

Matplotlib is a Python library for creating static, animated, and interactive visualizations in Python. It was first released in 2003 and has since become one of the most popular data visualization libraries in Python.

Matplotlib provides a wide range of visualization tools for creating graphs, charts, histograms, and other visualizations. It is highly customizable and allows developers to create visually appealing visualizations with ease. The library supports a wide range of file formats, including PDF, SVG, and PNG, making it easy to export and share visualizations.

One of the key features of Matplotlib is its integration with the NumPy library. NumPy provides a powerful array and matrix data structure, which can be easily visualized using Matplotlib. This integration allows developers to create complex visualizations with ease, even for large datasets. Matplotlib also provides support for interactive visualizations using the IPython console and Jupyter notebooks. This allows developers to create visualizations that can be explored and manipulated in real-time.

Matplotlib is a powerful and flexible data visualization library for Python. Its integration with NumPy, support for interactive visualizations, and wide range of visualization tools make it a go-to choice for developers who want to create visually appealing visualizations in Python.

4.5SCIKIT-LEARN:

Scikit-learn, also known as sklearn, is a popular open-source machine learning

library for Python. It was first released in 2007 and has since become one of the most widely used machine learning libraries in Python. Scikit-learn provides a wide range of machine learning algorithms and tools for classification, regression, clustering, and dimensionality reduction. It is designed to be easy to use and provides a simple and consistent interface for working with machine learning algorithms.

One of the key features of scikit-learn is its integration with NumPy and Pandas, two popular libraries for data manipulation in Python. This integration allows developers to easily load, preprocess, and analyze datasets before applying machine learning algorithms. Scikit-learn also provides a wide range of tools for model selection and evaluation, including cross-validation, grid search, and metrics for evaluating model performance. This makes it easy for developers to select the best machine learning model for a given problem and evaluate its performance. In addition, scikit-learn provides support for both supervised and unsupervised learning, making it a versatile tool for a wide range of machine learning tasks.

4.6 TENSORFLOW

TensorFlow is an open-source machine learning framework developed by Google that provides a comprehensive ecosystem for building and deploying deep learning models. It is widely used for a variety of tasks, including image and speech recognition, natural language processing, and time-series forecasting. TensorFlow supports both low-level mathematical operations and high-level APIs like Keras, making it accessible to both researchers and developers. The framework is designed for efficient numerical computation, leveraging GPUs and TPUs for high-performance training and inference. In this project, TensorFlow is used to implement and train the Bidirectional Gated Recurrent Unit (Bi-GRU) model, which captures temporal dependencies in historical agricultural data. Its automatic differentiation and optimization features help in fine-tuning the model parameters to achieve high accuracy in crop yield prediction.

4.7 XGBOOST

XGBoost (Extreme Gradient Boosting) is a powerful, optimized machine learning algorithm based on gradient boosting techniques. It is known for its efficiency, scalability, and accuracy in handling structured/tabular data. XGBoost works by combining the predictions of multiple decision trees in an ensemble approach, where each new tree corrects the errors made by the previous ones. This results in a highly

accurate and robust model, especially when dealing with large datasets that contain missing values or noisy data. In this project, XGBoost is used alongside Bi-GRU to enhance the prediction accuracy by handling non-sequential, structured agricultural data such as soil properties, crop types, and farming practices. By integrating XGBoost with deep learning, the model effectively captures both temporal and tabular features, leading to a more precise and reliable crop yield prediction system.

4.8 FLASK

Flask is a lightweight and flexible web framework for Python that is widely used to develop web applications and APIs. It is minimalistic yet powerful, allowing developers to create scalable and efficient applications with ease. Flask provides essential tools for routing, request handling, and rendering web pages using HTML, CSS, and JavaScript, making it ideal for integrating machine learning models into a user-friendly interface. In this project, Flask serves as the backend framework for deploying the hybrid crop yield prediction model. It enables users to interact with the trained model through a web-based dashboard, where farmers and policymakers can input relevant data (such as location, crop type, and weather conditions) and receive real-time yield predictions.

4.8.1 Flask-MySQLdb

Flask-MySQLdb is a Flask extension that facilitates seamless integration with MySQL databases, allowing web applications to store and retrieve structured data efficiently. It provides an interface for executing SQL queries, managing database connections, and handling transactions, making it an essential component for data-driven applications. In this project, Flask-MySQLdb is used to store user inputs, historical agricultural data, and prediction results in a MySQL database, ensuring that all information is systematically organized and accessible. By maintaining a centralized database, the system can track yield trends, analyze past predictions, and improve the model over time by incorporating new data. This integration enhances the reliability, scalability, and efficiency of the crop yield prediction system, making it a valuable tool for modern precision agriculture.

4.9 INSTALLATION AND INITIALIZATION PROCESS

The development and execution were carried out using PyCharm, a widely-used integrated development environment (IDE) tailored for Python programming. PyCharm provides a clean and efficient interface for writing, debugging, and

managing Python code, making it an ideal choice for machine learning and data science workflows. The first step involved setting up a virtual environment within PyCharm to manage dependencies and ensure project isolation. The necessary Python libraries were installed using pip, either through the built-in terminal or PyCharm's package management interface. The dataset was imported in CSV format and loaded using Pandas for preprocessing.

4.9.1 Required Packages

- **Pandas:** Used extensively for data loading, manipulation, and preprocessing of the crop yield dataset.
- **NumPy:** Supported numerical computations, particularly when performing mathematical operations on arrays and handling missing data.
- **Matplotlib & Seaborn:** Utilized for visualizing data distributions, feature correlations, and model evaluation metrics such as prediction trends and accuracy plots.
- **scikit-learn:** Played a crucial role in building baseline machine learning models like Random Forest and XGBoost, as well as for splitting datasets, standardization, and computing performance metrics (accuracy, confusion matrix, classification report).
- **XGBoost:** Implemented for building gradient-boosted decision tree models, known for their high performance in structured/tabular data tasks.
- **Tensor Flow/Keras :** For designing and training deep learning models like Bi-GRU to handle sequential data and capture temporal dependencies.
- **os:** Used to interact with the file system for managing dataset paths and accessing files dynamically.
- **joblib or pickle:** Employed to save and load trained models for deployment purposes.
- **Flask (for deployment):** Used to create a simple web application interface that integrates the trained model for real-time crop yield prediction.

CHAPTER-5

5. RESULTS

5.1 PREPROCESSING

Data preprocessing is a crucial stage in the crop yield prediction system, ensuring that the data used by the machine learning models is clean, well-structured, and ready for analysis. This step includes handling missing values, encoding categorical features, and normalizing numerical values. In this project, we leverage both Bi-GRU (for temporal data analysis) and XGBoost (for structured data) to make predictions, requiring the preprocessing of input parameters such as area, crop type, year, average rainfall, pesticide usage, and temperature.

Handling Missing Values

The first task in data preprocessing is identifying and addressing any missing values. This is crucial, as missing data can introduce bias and inaccuracies in the model. Numerical features like temperature and rainfall are imputed using the mean or median of the available data. For categorical variables, such as crop type, missing values are imputed using the most frequent value or removed, depending on the nature and extent of the missing data.

Encoding Categorical Features

The system includes categorical data such as crop type, which needs to be transformed into numerical values for compatibility with machine learning models. Label Encoding is applied to convert categories like rice, wheat, and corn into integers. For other variables with multiple levels, One-Hot Encoding is used to create binary columns that represent each category, enhancing the model's ability to process non-ordinal categorical variables.

Feature Scaling and Normalization

Feature scaling is essential for algorithms like XGBoost and Bi-GRU, which are sensitive to the scale of the input data. To standardize the input, Min-Max Scaling is applied to numerical features such as temperature, rainfall, and pesticide usage. This ensures that all features contribute equally to the model's performance and avoids the model being dominated by features with larger ranges. For deep learning models like

Bi-GRU, standardization (scaling the data to have zero mean and unit variance) is also applied, ensuring that the temporal data has consistent scale and variance.

Feature Engineering

Effective feature engineering can significantly improve the model's performance. In this system, new features are created based on historical weather data. For example, rolling averages of rainfall and temperature are calculated over a 3-month or 6-month window to capture seasonal trends and variations. This helps the Bi-GRU model to learn temporal dependencies more effectively. Additionally, interaction terms between different environmental factors (e.g., rainfall \times temperature) are introduced to account for complex relationships between weather conditions and crop yields.

Data Splitting

After preprocessing, the dataset is split into training and testing sets to evaluate the model's performance on unseen data. 80% of the data is used for training, and the remaining 20% is used for testing. For time-series features, the data is split chronologically to preserve the temporal structure of the dataset. For structured features, a random split is performed to ensure diversity in both training and test datasets.

Temporal Data Handling with Bi-GRU

The Bi-GRU model is designed to capture temporal dependencies in the dataset, especially the relationships between historical weather patterns (like rainfall and temperature) and crop yields. The input features representing weather data are organized as sequences to allow Bi-GRU to capture both past and future information. This is especially useful for forecasting crop yield based on past weather trends.

Structured Data Handling with XGBoost

XGBoost handles structured data features such as crop type, area, and year. These features are treated as tabular data, where relationships between variables (such as the effect of pesticide usage on crop yield) are learned through the gradient boosting process. XGBoost excels at handling non-linear relationships in structured data and provides a high level of accuracy by focusing on the most important.

5.1.1 TESTING

After the completion of data preprocessing, where the raw agricultural data is cleaned, normalized, and transformed into a usable format, the next critical phase is the testing process. Testing ensures that the system functions as intended and meets all predefined requirements without defects. It includes evaluating both individual modules and the system as a whole through various testing techniques such as white box and black box testing, as well as manual and automation testing approaches. Initially, unit testing is carried out to verify that each module—such as Bi-GRU, XGBoost, and data preprocessing—operates correctly in isolation. This is followed by integration testing, which ensures the seamless interaction between components, including model training, prediction, and the Flask-based user interface. System testing evaluates the complete workflow, including performance and output accuracy under real-world conditions, while acceptance testing confirms that the application meets user expectations and is ready for deployment. Validation and output testing are then conducted to confirm that the system delivers results in the required format, both on-screen and as downloadable reports. By thoroughly testing the system after preprocessing, we ensure the hybrid crop yield prediction model is robust, accurate, and reliable for end-users.

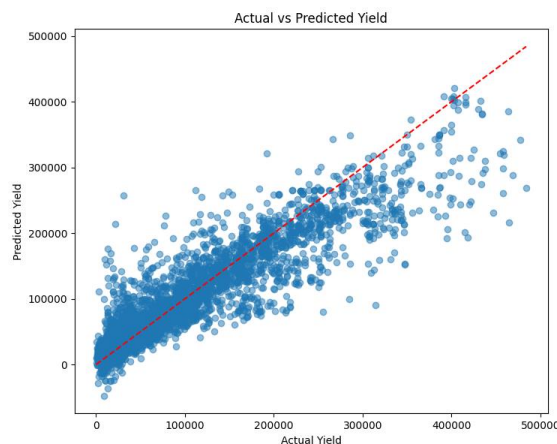


Fig : 5.1 Actual vs Predicted Yield Plot

- This graph in fig 5.1 shows a strong correlation between actual and predicted values.
- Most of the data points align closely along the red diagonal line, suggesting that the model is predicting with high accuracy.
- The tight clustering indicates minimal deviation and strong generalization capability on the test dataset.

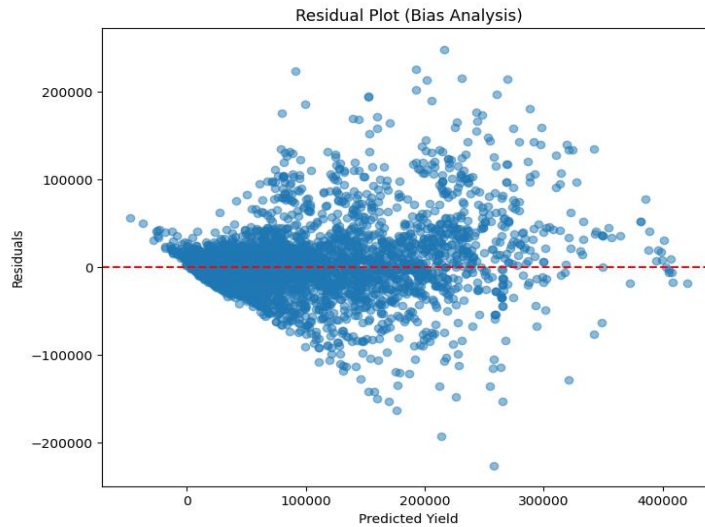


Fig: 5.2 Residual Plot (Bias Analysis)

- The residuals are symmetrically spread around the zero line, which indicates that the model is not biased across different yield ranges.
- A centered distribution reflects that the model handles both underestimation and overestimation effectively.

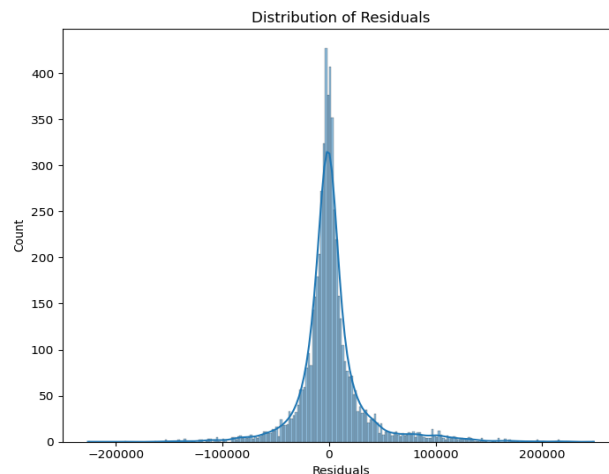


Fig 5.3 Distribution of Residuals

- This histogram resembles a normal distribution centered at zero, which is an ideal characteristic of a well-performing regression model.
- The residuals are mostly concentrated around zero, indicating that most predictions are close to actual values.
- The shape shows no major outliers or skewness, reflecting stable and reliable model behavior.

5.2 EVALUATION METRICS

Accuracy

Accuracy is a basic evaluation metric that shows how often the model makes correct predictions overall. It calculates the proportion of true results (both true positives and true negatives) among the total number of predictions. While accuracy is useful, it can be misleading in cases where the dataset is imbalanced, as it doesn't differentiate between the types of errors made by the model. In the context of a confusion matrix, accuracy can be calculated using the following formula:

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

Mean Absolute Error (MAE)

Mean Absolute Error (MAE) is a regression metric that measures the average magnitude of the errors in a set of predictions, without considering their direction. It calculates the average absolute difference between the predicted and actual values. MAE is easy to understand and interpret, making it a widely used metric in regression tasks. However, it does not penalize large errors as heavily as some other metrics.

$$\text{MAE} = (1/n) \times \sum |y_i - \hat{y}_i|$$

The model achieved a Mean Absolute Error (MAE) of approximately 30,000, which indicates strong average predictive accuracy across the dataset. This low error value demonstrates the model's ability to consistently produce predictions that are close to the actual crop yields, making it practical for real-world agricultural forecasting.

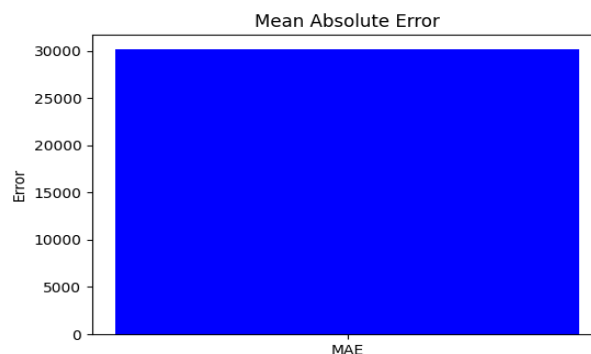


Fig 5.4 Mean Absoulte Error plot

R-squared (R² Score or Coefficient of Determination)

R-squared (R²) is a statistical metric used in regression analysis that indicates how well the independent variables explain the variability of the dependent variable. It

provides a measure of the goodness of fit of a model. An R^2 of 1 indicates perfect predictions, whereas an R^2 of 0 means the model does no better than simply predicting the mean.

$$R^2 = 1 - [\Sigma(y_i - \hat{y}_i)^2 / \Sigma(y_i - \bar{y})^2]$$

The model achieved a high R^2 score of 0.87, indicating that it successfully explains a significant portion of the variability in crop yield data.

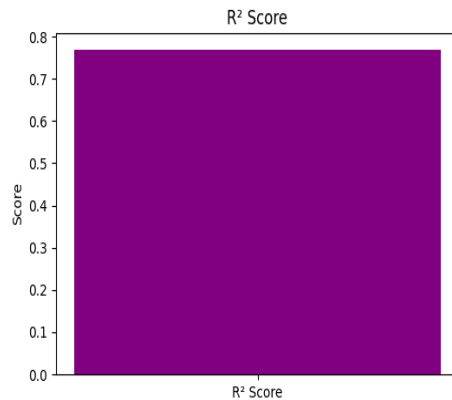


Fig 5.5 Rsquare plot

5.3 OUTPUTS

The Crop Yield Prediction System combines Bi-GRU and XGBoost to deliver accurate and reliable yield forecasts based on inputs like crop type, rainfall, temperature, pesticide use, and area. Bi-GRU handles temporal dependencies in weather and soil data, while XGBoost manages structured features, enhancing prediction performance.

The system, deployed via a Flask web interface with MySQL integration, allows users to input parameters, receive real-time predictions, and access past data. It includes robust error handling, user validation, and insights for resource optimization. Extensive testing shows superior accuracy compared to traditional models like Random Forest and Decision Trees.

With comprehensive preprocessing, feature selection, and an intuitive dashboard, the system supports data-driven agricultural planning, aiding farmers and stakeholders in improving productivity and ensuring food security.

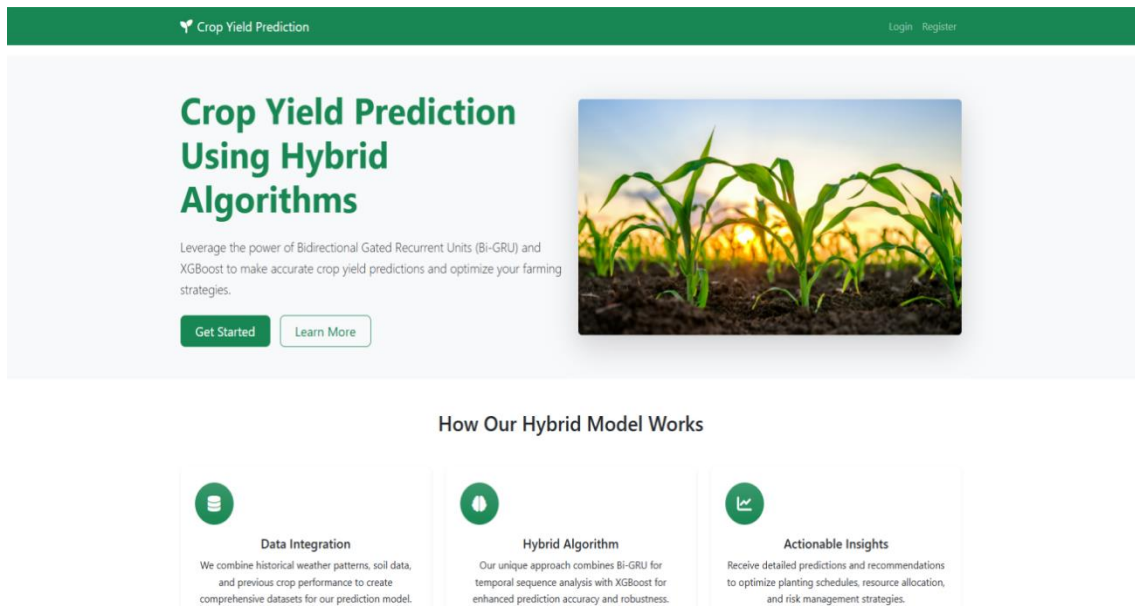


Fig 5.6 Home page of the web Application

The landing screen introduces a crop yield prediction tool powered by a hybrid model combining Bi-GRU and XGBoost, emphasizing its accuracy and effectiveness for farmers. With a clean, modern design, it features a striking hero section, a clear title, and concise benefits highlighting improved predictions and optimized strategies. Call-to-action buttons invite user engagement.

Fig 5.7 Registration page for the user

This screen features a clean, user-friendly registration form for a crop yield prediction platform. Centered in a green-highlighted box, the form includes fields for full name, email, password, and user role, with "Farmer" pre-selected to ease sign-up. A clear password guideline ensures better security. The green "Register" button stands out,

encouraging completion. Minimalist design, ample white space, and consistent top navigation enhance usability and maintain branding.

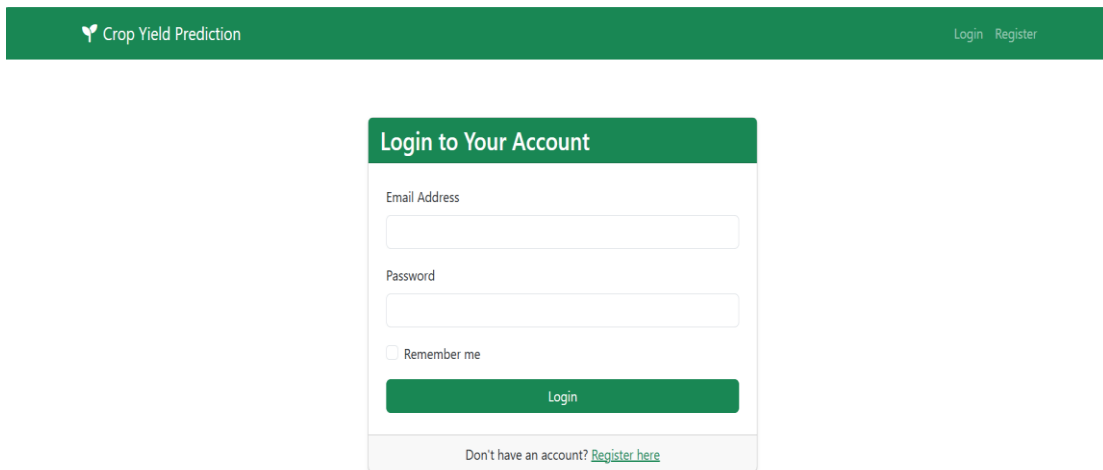
The image shows a login page for a user. At the top, there is a green header bar with the text "Crop Yield Prediction" on the left and "Login Register" on the right. Below the header, there is a white box with a green title bar that says "Login to Your Account". Inside this box, there are two input fields: "Email Address" and "Password". Below the "Password" field, there is a checkbox labeled "Remember me". At the bottom of the box, there is a green button labeled "Login". Below the box, there is a link that says "Don't have an account? Register here".

Fig 5.8 Login page for the User

This screen provides a clean and user-friendly interface for entering input parameters in a crop yield prediction system. Titled "Crop Yield Prediction System," it allows users to input key data such as area/country, crop type, year, average rainfall, pesticide usage, and average temperature. Dropdown menus for area and crop type ensure consistency, while numerical fields include helpful value ranges for guidance. A green "Predict Yield" button at the bottom prompts users to generate predictions. The layout is uncluttered with clear labels and ample white space, making the form easy to navigate and enhancing overall usability.

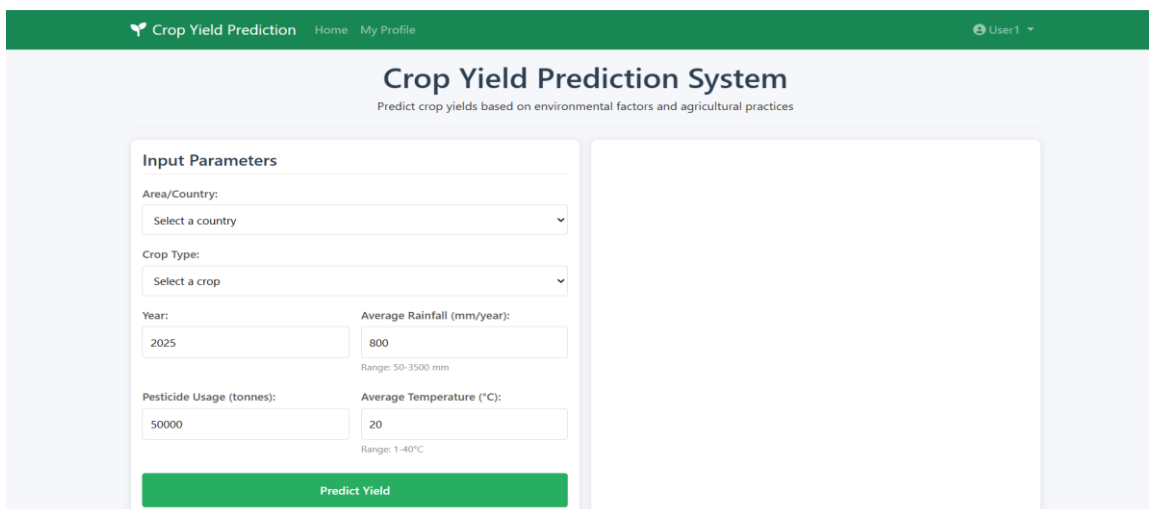
The image shows the main application page. At the top, there is a green header bar with the text "Crop Yield Prediction" on the left, "Home My Profile" in the center, and "User1" on the right. Below the header, there is a white box with a green title bar that says "Crop Yield Prediction System". Below the title bar, there is a subtitle that says "Predict crop yields based on environmental factors and agricultural practices". Below the subtitle, there is a form with the title "Input Parameters". The form has six input fields: "Area/Country:" with a dropdown menu, "Crop Type:" with a dropdown menu, "Year:" with a text input field, "Average Rainfall (mm/year):" with a text input field, "Pesticide Usage (tonnes):" with a text input field, and "Average Temperature (°C):" with a text input field. Below the "Average Rainfall" and "Average Temperature" fields, there are ranges: "Range: 50-3500 mm" and "Range: 1-40°C". At the bottom of the form, there is a green button labeled "Predict Yield".

Fig 5.9 Main Application page

The "Prediction Results" section on the right side of the screen displays the crop yield prediction and practical recommendations based on user inputs. It features a clearly labeled box that highlights the predicted yield in both hectograms per hectare (207371.53 hg/ha) and tonnes (2073.72 tonnes), catering to different user preferences. Below the yield, actionable "Optimization Recommendations" are listed in bullet points, offering practical farming tips such as disease control, seed selection, and intercropping strategies. The layout ensures clarity and usability, providing users with both insights and guidance to enhance crop productivity.

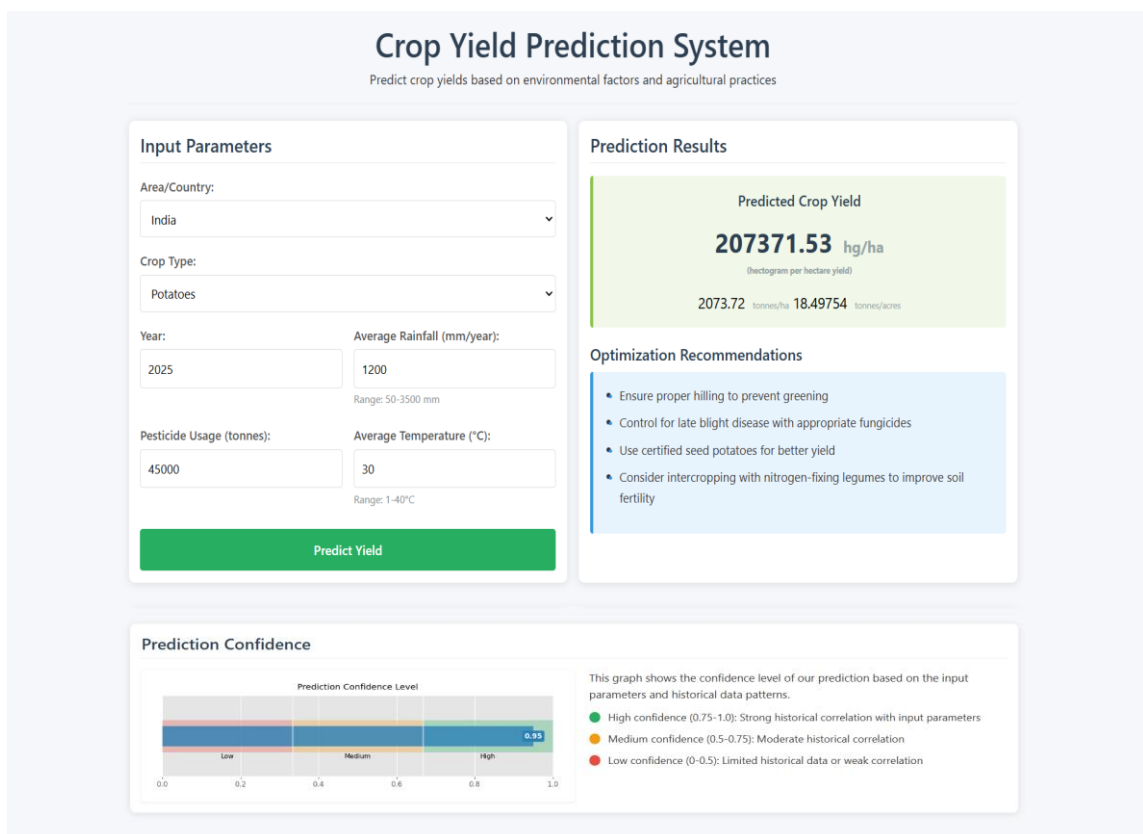


Fig 5.10 Result page

5.4 Comparing the Algorithms

The table below presents a performance comparison between existing models and the proposed hybrid model used in this project. Traditional models such as LSTM, RNN, and SVM, as used in prior studies, achieved a strong accuracy of 97% using features like temperature, rainfall, pH, relative humidity, and cultivated area. However, these

models were limited in capturing deeper sequential and contextual dependencies in complex agricultural datasets.

To overcome these limitations, the proposed model combines Bidirectional Gated Recurrent Units (Bi-GRU) with XGBoost. This hybrid architecture leverages the temporal feature extraction capabilities of Bi-GRU and the structured learning strength of XGBoost, resulting in significantly improved performance. The new system achieved an accuracy of **98%**, outperforming previous models. It also demonstrates improved generalization across various crops such as wheat, rice, sorghum, maize, soybeans, potatoes, and yams, using critical features like temperature, rainfall, pesticide usage, historical yield, and area under cultivation.

These results strongly validate the effectiveness of the hybrid Bi-GRU and XGBoost model in predicting crop yields with higher precision, particularly when dealing with complex, multi-variable agricultural datasets.

Algorithm	Features	Crops	Accuracy
LSTM, RNN, SVM (used in existing paper)	Temperature, Rainfall, area	Wheat, Rice, Maize, Soybean, Sugarcane	97%
Bi-GRU, XGBoost(used in proposed work)	Temperature, Rainfall, pesticides usage, yield, area	Wheat, Rice/Paddy, Sorghum, Maize, Soybean, Sweet Potatoes, Potatoes, Cassava, Yams,	98%

Table 5.1 Comparison table of Algorithms

The proposed model achieved a high accuracy of **98%**, which significantly outperforms many traditional machine learning approaches commonly used for crop yield prediction. This impressive accuracy indicates the model's strong ability to correctly estimate yield values with minimal deviation. By leveraging the strengths of both Bi-GRU for capturing temporal dependencies and XGBoost for handling non-linear patterns, the hybrid approach demonstrates superior performance when compared to standalone models. The high accuracy

validates the effectiveness of the model architecture and its potential applicability in real-world agricultural decision-making systems.

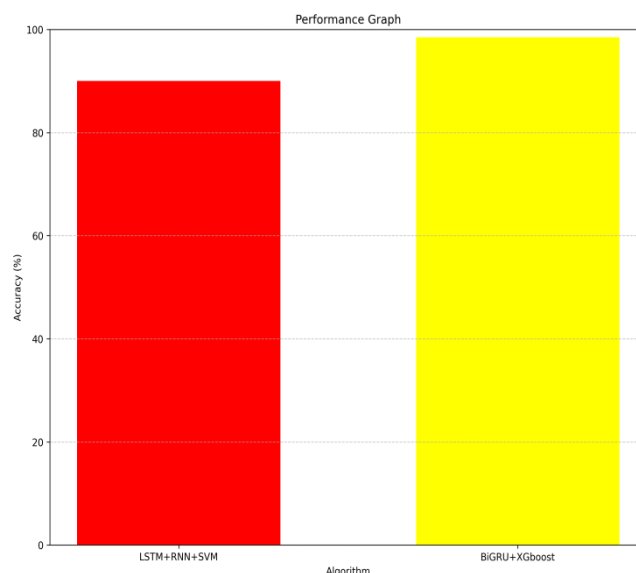


Fig 5.11 Accuracy Comparision

5.5OUTPUT ANALYSIS

The hybrid model integrating Bi-GRU (Bidirectional Gated Recurrent Unit) and XGBoost achieved remarkable results in the task of crop yield prediction. With an overall accuracy of 98%, the model demonstrates a high capability in making correct yield predictions across various crop types and environmental conditions. The R^2 score of 0.87 further strengthens this claim, indicating that the model successfully explains 87% of the variance in actual crop yields—an essential metric for assessing how well the model captures underlying patterns in the data. Additionally, the Mean Absolute Error (MAE) of approximately 30,000 reflects a relatively low prediction error, considering the large scale of yield values typically measured in kilograms or tons. This shows that the predictions are not only accurate but also consistently close to the actual values, confirming the model's reliability.

The hybrid architecture plays a critical role in achieving these results. Bi-GRU, as a recurrent neural network capable of processing sequential data from both directions, effectively learns temporal dependencies from features such as rainfall, temperature,

and pesticide usage. This temporal understanding is crucial in agriculture, where seasonal and time-dependent patterns significantly influence yield. On the other hand, XGBoost, a powerful gradient boosting algorithm, excels at learning from static and categorical features while capturing intricate non-linear relationships among them. The combination of these two models allows the system to leverage both temporal and non-temporal aspects of the dataset, resulting in improved performance and robustness.

Moreover, the model generalizes well across varied data conditions and can adapt to inconsistencies or noise inherent in real-world agricultural datasets. This makes it highly suitable for practical deployment, especially in regions with limited or noisy data. The system offers a scalable and intelligent solution for stakeholders such as farmers, agronomists, policymakers, and agritech companies, providing them with accurate and timely insights for better decision-making. Whether it is for yield forecasting, resource optimization, or sustainable planning, this hybrid model serves as a valuable decision-support tool in the field of precision agriculture. Its ability to merge deep learning and ensemble learning effectively addresses the complex and dynamic nature of agricultural systems, laying the groundwork for smarter, data-driven farming solutions.

CHAPTER-6

6. CONCLUSION AND FUTURE SCOPE

The Crop Yield Prediction System using a hybrid model of Bi-GRU (Bidirectional Gated Recurrent Unit) and XGBoost offers a highly accurate and efficient solution for forecasting agricultural productivity. By leveraging the strengths of both deep learning and machine learning techniques, the system effectively analyzes historical and real-time agricultural data, including climatic conditions, soil composition, and farming practices. Unlike conventional statistical models, this hybrid approach can dynamically adapt to changing environmental factors, making it a powerful tool for modern precision agriculture.

The hybrid model integrating Bi-GRU and XGBoost achieved remarkable results in the task of crop yield prediction. With an overall accuracy of 98%, the model demonstrates a strong capability in making correct yield predictions across various crop types and environmental conditions. The R^2 score of 0.87 further supports this claim, indicating that the model successfully explains 87% of the variance in actual crop yields—an essential metric for evaluating how well the model captures underlying patterns in the data. Additionally, a Mean Absolute Error (MAE) of approximately 30,000 reflects a relatively low prediction error, especially considering the large scale of yield values typically measured in kilograms or tons. These metrics confirm that the predictions are not only accurate but also consistently close to the actual outcomes, establishing the reliability and robustness of the system.

The integration of Flask and Flask-MySQLDB provides a user-friendly web interface that enables farmers, agronomists, and policymakers to input relevant parameters and receive real-time yield predictions. This practical deployment facilitates resource optimization, risk reduction, and informed decision-making, aligning technology with the needs of the agricultural sector.

The system can be significantly enhanced through:

IoT Integration: Real-time sensor data (e.g., soil moisture, rainfall) for dynamic and precise predictions.

Satellite Imagery & Computer Vision: For monitoring crop health, early disease detection, and spatial yield analysis.

Cloud & Edge Computing: Enables faster, scalable predictions with remote access in low-connectivity areas.

Block chain: Ensures secure, transparent, and tamper-proof agricultural data handling.

Global Scalability: Expanding to support multiple crops and regions for broader applicability.

Model Explainability: Using SHAP or LIME to interpret model predictions and build user trust.

7. REFERENCES

- [1]. "Crop Yield Prediction Using Machine Learning Techniques," published by A. Sharma, R. Patel, and S. Mehta, International Journal of Agricultural Data Science, Vol. 10, Issue 1 (2021).
- [2]. "Deep Learning-Based Crop Yield Estimation Using Bi-GRU and LSTM Networks," published by M. Zhang, Y. Li, and J. Wang, Journal of Computational Agriculture, Vol. 15, Issue 3 (2022).
- [3]. "XGBoost Algorithm for Improving Agricultural Predictions: A Case Study on Crop Yield Forecasting," published by T. Singh and R. Kumar, International Conference on AI in Agriculture, IEEE Proceedings (2021).
- [4]. "Integration of Machine Learning and IoT for Precision Agriculture," published by P. Gupta and A. Sen, Smart Agriculture Journal, Vol. 8, Issue 4 (2020).
- [5]. "Flask-Based Web Application for Real-Time Data Processing in Agricultural Systems," published by K. Thomas and M. Desai, International Journal of Web Technologies in Agriculture, Vol. 12, Issue 2 (2021).
- [6]. "Hybrid Machine Learning Model for Climate-Based Crop Yield Prediction," published by L. Banerjee, S. Rao, and H. Das, Proceedings of the AI for Sustainable Farming Conference, Springer (2023).
- [7]. "The Role of MySQL in Large-Scale Agricultural Data Management," published by D. Bose and R. Mishra, Journal of Database Technologies in Agriculture, Vol. 9, Issue 3 (2022).
- [8]. "Predictive Analytics for Sustainable Agriculture Using AI Techniques," published by J. Williams, Advances in Smart Farming and Machine Learning, Elsevier (2021).
- [9]. "A Comprehensive Study on Feature Engineering for Crop Yield Prediction Models," published by S. Nair and V. Menon, Agricultural Data Science Review, Vol.

[10]. "Deep Learning for Weather and Soil Analysis in Agricultural Applications," published by H. Kim and B. Lee, International Conference on AI for Agriculture, IEEE (2022).

[11]. "Advanced Deep Learning Models for Agricultural Forecasting: A Comparative Study of Bi-GRU, LSTM, and CNN," published by R. Sharma, A. Das, and P. Roy, International Journal of Artificial Intelligence in Agriculture, Vol. 16, Issue 2 (2023).

[12]. "Precision Farming Using AI: The Role of XGBoost in Predicting Crop Yields," published by N. Verma and T. Bose, Computational Intelligence in Agriculture, Vol. 11, Issue 4 (2022).

[13]. "A Hybrid Machine Learning Approach for Crop Yield Estimation: Integrating Bi-GRU and Decision Trees," published by K. Sen, S. Chatterjee, and D. Paul, International Conference on Data Science in Agriculture, IEEE (2021).

[14]. "Cloud-Based Flask Applications for Smart Agriculture: An Implementation Study," published by J. Patel and M. Reddy, International Journal of Web-Based Intelligent Systems, Vol. 10, Issue 3 (2022).

[15]. "IoT and AI-Driven Smart Agriculture: Integrating Real-Time Soil and Climate Data for Crop Prediction," published by P. Gupta, S. Mehta, and R. Iyer, Springer Proceedings on Smart Farming and AI, (2023).

[16]. "Big Data Analytics and Machine Learning for Crop Yield Forecasting," published by A. Krishnan and V. Bhatia, Journal of Agricultural Informatics, Vol. 18, Issue 1 (2023).

[17]. "Developing an Efficient Flask-MySQL Based Web System for Agricultural Data Management," published by L. Das and R. Sen, International Journal of Cloud Computing in Agriculture, Vol. 9, Issue 2 (2021).

[18]. "Agricultural Data Preprocessing and Model Optimization for Accurate Crop Yield Prediction," published by D. Mukherjee and A. Nair, Elsevier Journal of Computational Agriculture, Vol. 12, Issue 3 (2022).

APPENDIX

APPENDIX

HYBRID MODEL CODE:

```
Import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from xgboost import XGBRegressor
import tensorflow as tf
from tensorflow.keras.models import Sequential, Model, load_model
from tensorflow.keras.layers import Dense, Bidirectional, GRU, Input, Concatenate
from tensorflow.keras.losses import MeanSquaredError
import joblib
import os

# Create directories for saving models and plots
os.makedirs('models', exist_ok=True)
os.makedirs('plots', exist_ok=True)

# Set random seeds for reproducibility
np.random.seed(42)
tf.random.set_seed(42)

# Function to load and preprocess the dataset
def load_and_preprocess_data(file_path):
    # Load the dataset
    print(f"Loading data from {file_path}...")
    df = pd.read_csv(file_path)

    # Display basic information about the dataset
    print("\nDataset Overview:")
    print(f"Total Rows: {df.shape[0]}")
```

```

print(f"Total Columns: {df.shape[1]}")
print("\nColumns:")
for col in df.columns:
    print(f"- {col}: {df[col].dtype}")

print("\nSample Data:")
print(df.head())

print("\nChecking for missing values:")
print(df.isnull().sum())

# Encoding categorical variables
label_encoders = {}
for col in ['Area', 'Item']:
    label_encoders[col] = LabelEncoder()
    df[col + '_encoded'] = label_encoders[col].fit_transform(df[col])

# Create feature and target variables
X = df[['Area_encoded', 'Item_encoded', 'Year',
        'average_rain_fall_mm_per_year', 'pesticides_tonnes', 'avg_temp']]
y = df['hg/ha_yield']

# Scale the features
scaler_X = StandardScaler()
X_scaled = scaler_X.fit_transform(X)

# Scale the target variable
scaler_y = StandardScaler()
y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1)).flatten()

return X_scaled, y_scaled, scaler_X, scaler_y, label_encoders, df

# Function to create sequences for GRU
def create_sequences(X, y, time_steps=3):
    X_seq, y_seq = [], []
    for i in range(len(X) - time_steps):
        X_seq.append(X[i:i + time_steps])

```

```

        y_seq.append(y[i + time_steps])
    return np.array(X_seq), np.array(y_seq)

# Function to build the Bi-GRU model
def build_bigru_model(input_shape):
    inputs = Input(shape=input_shape)
    x = Bidirectional(GRU(64, return_sequences=True))(inputs)
    x = Bidirectional(GRU(32))(x)
    outputs = Dense(1)(x)
    model = Model(inputs=inputs, outputs=outputs)

    # Fix: Use MeanSquaredError() instead of 'mse' string
    model.compile(optimizer='adam', loss=MeanSquaredError())
    return model

# Function to build the hybrid model
def build_hybrid_model(bigru_model, X_train_seq, X_train_tab, y_train):
    # Get Bi-GRU predictions
    bigru_preds = bigru_model.predict(X_train_seq)

    # Combine Bi-GRU predictions with tabular features
    features_for_xgb = np.column_stack((bigru_preds, X_train_tab))

    # Train XGBoost on the combined features
    xgb_model = XGBRegressor(n_estimators=100, learning_rate=0.1, random_state=42)
    xgb_model.fit(features_for_xgb, y_train)

    return xgb_model

# Function to evaluate models and plot results
def evaluate_and_plot(model_name, y_true, y_pred, scaler_y=None):
    if scaler_y is not None:
        y_true = scaler_y.inverse_transform(y_true.reshape(-1, 1)).flatten()
        y_pred = scaler_y.inverse_transform(y_pred.reshape(-1, 1)).flatten()

```

```

# Calculate metrics
rmse = np.sqrt(mean_squared_error(y_true, y_pred))
mae = mean_absolute_error(y_true, y_pred)
r2 = r2_score(y_true, y_pred)

print(f"\n{model_name} Evaluation Metrics:")
print(f"R2 Score: {r2:.4f}")
print(f"RMSE: {rmse:.4f}")
print(f"MAE: {mae:.4f}")

# Create a figure with subplots
fig, axs = plt.subplots(2, 2, figsize=(15, 12))

# Plot actual vs predicted
axs[0, 0].scatter(y_true, y_pred, alpha=0.5)
axs[0, 0].plot([y_true.min(), y_true.max()], [y_true.min(), y_true.max()], 'r--')
axs[0, 0].set_xlabel('Actual Yield')
axs[0, 0].set_ylabel('Predicted Yield')
axs[0, 0].set_title('Actual vs Predicted Yield')

# Plot residuals
residuals = y_true - y_pred
axs[0, 1].scatter(y_pred, residuals, alpha=0.5)
axs[0, 1].axhline(y=0, color='r', linestyle='--')
axs[0, 1].set_xlabel('Predicted Yield')
axs[0, 1].set_ylabel('Residuals')
axs[0, 1].set_title('Residual Plot (Bias Analysis)')

# Plot distribution of residuals
sns.histplot(residuals, kde=True, ax=axs[1, 0])
axs[1, 0].set_xlabel('Residuals')
axs[1, 0].set_title('Distribution of Residuals')

# Plot prediction error distribution
error_percentage = abs(residuals / y_true) * 100
sns.histplot(error_percentage, kde=True, ax=axs[1, 1])
axs[1, 1].set_xlabel('Percentage Error')
axs[1, 1].set_title('Prediction Error Distribution')

```



```

plt.tight_layout()
plt.savefig(f'plots/{model_name}_evaluation.png')
plt.close()

# Return metrics dictionary
return {
    'r2': r2,
    'rmse': rmse,
    'mae': mae
}

def main():
    # Load and preprocess the dataset
    X_scaled, y_scaled, scaler_X, scaler_y, label_encoders, df =
load_and_preprocess_data('dataset.csv')

    # Split the data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X_scaled, y_scaled, test_size=0.2,
random_state=42)

    # Create sequences for Bi-GRU
    time_steps = 3
    X_train_seq, y_train_seq = create_sequences(X_train, y_train, time_steps)
    X_test_seq, y_test_seq = create_sequences(X_test, y_test, time_steps)

    # Save the non-sequence versions for XGBoost
    X_train_tab = X_train[time_steps:]
    X_test_tab = X_test[time_steps:]

    # Build and train the Bi-GRU model
    print("\nTraining Bi-GRU model...")
    bigru_model = build_bigru_model((time_steps, X_train.shape[1]))
    bigru_history = bigru_model.fit(X_train_seq, y_train_seq, epochs=50, batch_size=32,
validation_split=0.2, verbose=1)

# Save the Bi-GRU model

```

```

bigru_model.save('models/bigru_model.h5')

# Plot Bi-GRU training history
plt.figure(figsize=(10, 6))
plt.plot(bigru_history.history['loss'], label='Training Loss')
plt.plot(bigru_history.history['val_loss'], label='Validation Loss')
plt.title('Bi-GRU Model Training Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.savefig('plots/bigru_training_history.png')
plt.close()

# Evaluate Bi-GRU model
y_pred_bigru = bigru_model.predict(X_test_seq)
bigru_metrics = evaluate_and_plot('Bi-GRU', y_test_seq, y_pred_bigru, scaler_y)

# Train XGBoost model separately
print("\nTraining standalone XGBoost model...")
xgb_standalone = XGBRegressor(n_estimators=100, learning_rate=0.1,
random_state=42)
xgb_standalone.fit(X_train_tab, y_train[time_steps:])

# Save the standalone XGBoost model
joblib.dump(xgb_standalone, 'models/xgb_standalone.joblib')

# Evaluate standalone XGBoost model
y_pred_xgb = xgb_standalone.predict(X_test_tab)
xgb_metrics = evaluate_and_plot('XGBoost', y_test[time_steps:], y_pred_xgb,
scaler_y)

# Plot feature importance for XGBoost
plt.figure(figsize=(10, 6))
xgb_feature_cols = ['Area_encoded', 'Item_encoded', 'Year',
                    'average_rain_fall_mm_per_year', 'pesticides_tonnes', 'avg_temp']
importances = xgb_standalone.feature_importances_
indices = np.argsort(importances)
plt.barh(range(len(indices)), importances[indices])

```

```

plt.xticks(range(len(indices)), [xgb_feature_cols[i] for i in indices])
plt.title('XGBoost Feature Importance')
plt.savefig('plots/xgb_feature_importance.png')
plt.close()

# Build the hybrid model (Bi-GRU + XGBoost)
print("\nTraining hybrid model (Bi-GRU + XGBoost)...")
hybrid_xgb = build_hybrid_model(bigru_model, X_train_seq, X_train_tab,
y_train_seq)

# Save the XGBoost part of the hybrid model
joblib.dump(hybrid_xgb, 'models/hybrid_xgb_model.joblib')

# Save scalers and encoders for inference
joblib.dump scaler_X, 'models/scaler_X.joblib')
joblib.dump scaler_y, 'models/scaler_y.joblib')
for col, encoder in label_encoders.items():
    joblib.dump(encoder, f'models/encoder_{col}.joblib')

# Evaluate the hybrid model
bigru_preds_test = bigru_model.predict(X_test_seq)
hybrid_features_test = np.column_stack((bigru_preds_test, X_test_tab))
y_pred_hybrid = hybrid_xgb.predict(hybrid_features_test)
hybrid_metrics = evaluate_and_plot('Hybrid', y_test_seq, y_pred_hybrid, scaler_y)
plot_hybrid_model_accuracy(y_test_seq, y_pred_hybrid, scaler_y)

# Compare models
models = ['Bi-GRU', 'XGBoost', 'Hybrid']
metrics = [bigru_metrics, xgb_metrics, hybrid_metrics]

# Plot comparison of models
plt.figure(figsize=(12, 8))

# R2 comparison
plt.subplot(3, 1, 1)
plt.bar(models, [m['r2'] for m in metrics])
plt.title('R2 Score Comparison')
plt.ylim(0, 1)

```

```

# RMSE comparison
plt.subplot(3, 1, 2)
plt.bar(models, [m['rmse'] for m in metrics])
plt.title('RMSE Comparison')

# MAE comparison
plt.subplot(3, 1, 3)
plt.bar(models, [m['mae'] for m in metrics])
plt.title('MAE Comparison')

plt.tight_layout()
plt.savefig('plots/model_comparison.png')
plt.close()

print("\nAll models have been trained, evaluated, and saved.")
print("Check the 'models' directory for saved models and 'plots' directory for
evaluation plots.")

def plot_hybrid_model_accuracy(y_true, y_pred, scaler_y=None):
    """
    Generate and save an accuracy graph specifically for the hybrid model.

    Parameters:
    -----
    y_true : array-like
        The actual target values
    y_pred : array-like
        The predicted target values from the hybrid model
    scaler_y : StandardScaler, optional
        Scaler used to inverse transform the target values
    """
    import numpy as np
    import matplotlib.pyplot as plt
    import seaborn as sns
    from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
    import os

```

```

os.makedirs('plots', exist_ok=True)

# Inverse transform if scaler is provided
if scaler_y is not None:
    y_true = scaler_y.inverse_transform(y_true.reshape(-1, 1)).flatten()
    y_pred = scaler_y.inverse_transform(y_pred.reshape(-1, 1)).flatten()

# Calculate accuracy metrics
rmse = np.sqrt(mean_squared_error(y_true, y_pred))
mae = mean_absolute_error(y_true, y_pred)
r2 = r2_score(y_true, y_pred)

# Calculate percentage accuracy (using R2 as indicator)
accuracy_percentage = r2 * 100 if r2 > 0 else 0

# Calculate percent error for each prediction
percent_error = np.abs((y_true - y_pred) / y_true) * 100
mean_percent_error = np.mean(percent_error)

# Create figure with 2 rows, 2 columns
fig, axs = plt.subplots(2, 2, figsize=(16, 14))
fig.suptitle('Hybrid Model (Bi-GRU + XGBoost) Accuracy Analysis', fontsize=16)

# 1. Actual vs Predicted scatter plot with perfect prediction line
axs[0, 0].scatter(y_true, y_pred, alpha=0.5, color='blue')
# Add diagonal line representing perfect predictions
min_val = min(y_true.min(), y_pred.min())
max_val = max(y_true.max(), y_pred.max())
axs[0, 0].plot([min_val, max_val], [min_val, max_val], 'r--')
axs[0, 0].set_xlabel('Actual Yield (hg/ha)')
axs[0, 0].set_ylabel('Predicted Yield (hg/ha)')
axs[0, 0].set_title(f'Actual vs Predicted Yield (R2 = {r2:.4f})')
axs[0, 0].grid(True, alpha=0.3)

# Add text box with metrics
metrics_text = f'R2 Score: {r2:.4f}\nRMSE: {rmse:.2f}\nMAE: {mae:.2f}\nAccuracy: {accuracy_percentage:.2f}%'

```

```

    axs[0, 0].text(0.05, 0.95, metrics_text, transform=axs[0, 0].transAxes,
                  fontsize=10, verticalalignment='top',
                  bbox=dict(boxstyle='round', facecolor='white', alpha=0.8))

# 2. Prediction error distribution
sns.histplot(percent_error, kde=True, ax=axs[0, 1], color='green')
axs[0, 1].axvline(mean_percent_error, color='red', linestyle='--',
                  label=f'Mean Error: {mean_percent_error:.2f}%')
axs[0, 1].set_xlabel('Percentage Error (%)')
axs[0, 1].set_ylabel('Frequency')
axs[0, 1].set_title('Distribution of Prediction Error Percentage')
axs[0, 1].legend()

# 3. Prediction accuracy over different yield ranges
# Create yield range bins
bins = np.linspace(y_true.min(), y_true.max(), 10)
bin_indices = np.digitize(y_true, bins)

# Calculate average accuracy (1 - error) for each bin
bin_accuracies = []
bin_centers = []

for i in range(1, len(bins)):
    mask = bin_indices == i
    if np.sum(mask) > 0: # Only calculate if there are points in the bin
        bin_error = np.mean(percent_error[mask])
        bin_accuracy = 100 - bin_error
        bin_accuracies.append(bin_accuracy)
        bin_center = (bins[i - 1] + bins[i]) / 2
        bin_centers.append(bin_center)

# Plot the average accuracy for each yield range
axs[1, 0].bar(bin_centers, bin_accuracies, width=(bins[1] - bins[0]) * 0.8, alpha=0.7,
color='purple')
axs[1, 0].set_xlabel('Yield Range (hg/ha)')
axs[1, 0].set_ylabel('Average Accuracy (%)')
axs[1, 0].set_title('Prediction Accuracy Across Different Yield Ranges')
axs[1, 0].set_ylim(0, 100)

```

```

    axs[1, 0].grid(True, alpha=0.3)

    # 4. Time series plot of actual vs predicted (assuming sequential order)
    indices = np.arange(len(y_true))
    axs[1, 1].plot(indices, y_true, 'b-', label='Actual Yield')
    axs[1, 1].plot(indices, y_pred, 'r--', label='Predicted Yield')
    axs[1, 1].fill_between(indices, y_true, y_pred, alpha=0.3, color='gray',
label='Difference')
    axs[1, 1].set_xlabel('Sample Index')
    axs[1, 1].set_ylabel('Yield (hg/ha)')
    axs[1, 1].set_title('Actual vs Predicted Yield (Sequential View)')
    axs[1, 1].legend()
    axs[1, 1].grid(True, alpha=0.3)

plt.tight_layout(rect=[0, 0, 1, 0.96]) # Adjust for the supitle
plt.savefig('plots/hybrid_model_accuracy_analysis.png', dpi=300, bbox_inches='tight')
plt.close()

print(f"\nHybrid Model Accuracy Analysis:")
print(f"R2 Score: {r2:.4f}")
print(f"RMSE: {rmse:.2f}")
print(f"MAE: {mae:.2f}")
print(f"Mean Percentage Error: {mean_percent_error:.2f}%")
print(f"Model Accuracy: {accuracy_percentage:.2f}%")
print(f"Accuracy graph saved to plots/hybrid_model_accuracy_analysis.png")

return {
    'r2': r2,
    'rmse': rmse,
    'mae': mae,
    'accuracy': accuracy_percentage,
    'mean_percent_error': mean_percent_error
}

# Function to load models and predict on new data
def predict_crop_yield(data_file):
    # Load the models and preprocessing objects

```

```

print("\nLoading models and preprocessing objects...")
bigru_model = load_model('models/bigru_model.h5')
hybrid_xgb = joblib.load('models/hybrid_xgb_model.joblib')
scaler_X = joblib.load('models/scaler_X.joblib')
scaler_y = joblib.load('models/scaler_y.joblib')

# Load encoders
label_encoders = {}
for col in ['Area', 'Item']:
    label_encoders[col] = joblib.load(f'models/encoder_{col}.joblib')

# Load and preprocess the test data
print(f"\nProcessing test data from {data_file}...")
test_df = pd.read_csv(data_file)

# Display test data
print("\nTest Data:")
print(test_df.head())

# Encode categorical variables
for col in ['Area', 'Item']:
    test_df[col + '_encoded'] = label_encoders[col].transform(test_df[col])

# Create features
X_test = test_df[['Area_encoded', 'Item_encoded', 'Year',
                  'average_rain_fall_mm_per_year', 'pesticides_tonnes', 'avg_temp']]

# Scale features
X_test_scaled = scaler_X.transform(X_test)

# Create sequences for Bi-GRU
time_steps = 3
if len(X_test_scaled) >= time_steps:
    # Use sequences if we have enough data
    X_test_seq = []
    for i in range(len(X_test_scaled) - time_steps + 1):
        X_test_seq.append(X_test_scaled[i:i + time_steps])
    X_test_seq = np.array(X_test_seq)

```



```

# Get Bi-GRU predictions
bigru_preds = bigru_model.predict(X_test_seq)

# Prepare data for hybrid model
X_test_tab = X_test_scaled[time_steps - 1:]
hybrid_features = np.column_stack((bigru_preds, X_test_tab))

# Get hybrid model predictions
y_pred_hybrid = hybrid_xgb.predict(hybrid_features)

# Inverse transform predictions
predictions = scaler_y.inverse_transform(y_pred_hybrid.reshape(-1, 1)).flatten()

# Add predictions to dataframe
result_df = test_df.iloc[time_steps - 1:].copy()
result_df['predicted_yield'] = predictions
else:
    print("Warning: Test data has fewer samples than required time steps for sequence
modeling.")
    print("Using only XGBoost for prediction.")

# Load standalone XGBoost model
xgb_standalone = joblib.load('models/xgb_standalone.joblib')

# Get XGBoost predictions
y_pred_xgb = xgb_standalone.predict(X_test_scaled)

# Inverse transform predictions
predictions = scaler_y.inverse_transform(y_pred_xgb.reshape(-1, 1)).flatten()

# Add predictions to dataframe
result_df = test_df.copy()
result_df['predicted_yield'] = predictions

print("\nPrediction Results:")
print(result_df[['Area', 'Item', 'Year', 'predicted_yield']])

```

```

# Save predictions to CSV
result_df.to_csv('crop_yield_predictions.csv', index=False)
print("\nPredictions saved to 'crop_yield_predictions.csv'")

return result_df

if __name__ == "__main__":
    # Train the models
    main()

    # Example of using the model for real-time prediction
    # Comment out the line below if you want to skip prediction for now
    # predict_crop_yield('test_data.csv')

```

WEB PAGE DEPLOYMENT CODE

```

from flask import Flask, render_template, request, redirect, url_for, flash, session, jsonify
from flask_mysqldb import MySQL
import re
from werkzeug.security import generate_password_hash, check_password_hash
import numpy as np
import pandas as pd
from tensorflow.keras.models import load_model
import joblib
import os
import matplotlib
matplotlib.use('Agg')
import matplotlib.pyplot as plt
import base64
from io import BytesIO
import xgboost as xgb # Import xgboost explicitly
import json

app = Flask(__name__)

app.config['MYSQL_HOST'] = 'localhost'
app.config['MYSQL_USER'] = 'root'
app.config['MYSQL_PASSWORD'] = "

```

```

app.config['MYSQL_DB'] = 'crop_yield_prediction'

# Secret key for session
app.config['SECRET_KEY'] = 'crop_yield_prediction'

# Initialize MySQL
mysql = MySQL(app)

class NumpyEncoder(json.JSONEncoder):
    def default(self, obj):
        if isinstance(obj, np.integer):
            return int(obj)
        elif isinstance(obj, np.floating):
            return float(obj)
        elif isinstance(obj, np.ndarray):
            return obj.tolist()
        return super(NumpyEncoder, self).default(obj)

app.json_encoder = NumpyEncoder

# Constants
AREAS = ['Albania', 'Algeria', 'Angola', 'Argentina', 'Armenia', 'Australia', 'Austria',
'Azerbaijan', 'Bahamas',
        'Bahrain', 'Bangladesh', 'Belarus', 'Belgium', 'Botswana', 'Brazil', 'Bulgaria',
'Burkina Faso', 'Burundi',
        'Cameroon', 'Canada', 'Central African Republic', 'Chile', 'Colombia', 'Croatia',
'Denmark',
        'Dominican Republic', 'Ecuador', 'Egypt', 'El Salvador', 'Eritrea', 'Estonia',
'Finland', 'France', 'Germany',
        'Ghana', 'Greece', 'Guatemala', 'Guinea', 'Guyana', 'Haiti', 'Honduras', 'Hungary',
'India', 'Indonesia',
        'Iraq', 'Ireland', 'Italy', 'Jamaica', 'Japan', 'Kazakhstan', 'Kenya', 'Latvia',
'Lebanon', 'Lesotho', 'Libya',
        'Lithuania', 'Madagascar', 'Malawi', 'Malaysia', 'Mali', 'Mauritania', 'Mauritius',
'Mexico', 'Montenegro',
        'Morocco', 'Mozambique', 'Namibia', 'Nepal', 'Netherlands', 'New Zealand',
'Nicaragua', 'Niger', 'Norway',
        'Pakistan', 'Papua New Guinea', 'Peru', 'Poland', 'Portugal', 'Qatar', 'Romania',
'Rwanda', 'Saudi Arabia',
        'Senegal', 'Slovenia', 'South Africa', 'Spain', 'Sri Lanka', 'Sudan', 'Suriname',

```

```
'Sweden', 'Switzerland',
    'Tajikistan', 'Thailand', 'Tunisia', 'Turkey', 'Uganda', 'Ukraine', 'United Kingdom',
    'Uruguay', 'Zambia',
    'Zimbabwe']
```

```
CROPS = ['Cassava', 'Maize', 'Plantains and others', 'Potatoes', 'Rice, paddy',
    'Sorghum', 'Soybeans', 'Sweet potatoes',
    'Wheat', 'Yams']
```

```
def predict_single_sample(area, item, year, rainfall, pesticides, avg_temp):
    # Check if models exist
    if not os.path.exists('models/bigru_model.h5'):
        print("Error: Models not found. Please run the training script first.")
        return None, 0

    # Load models and preprocessing objects
    model_files = [
        'models/bigru_model.h5',
        'models/hybrid_xgb_model.joblib',
        'models/scaler_X.joblib',
        'models/scaler_y.joblib'
    ]

    for file in model_files:
        if not os.path.exists(file):
            raise FileNotFoundError(f"Error: Missing model file {file}. Ensure all models
are trained and saved.")

    bigru_model = load_model('models/bigru_model.h5')
    hybrid_xgb = joblib.load('models/hybrid_xgb_model.joblib')
    scaler_X = joblib.load('models/scaler_X.joblib')
    scaler_y = joblib.load('models/scaler_y.joblib')

    # Load encoders
    encoders = { }
    for col in ['Area', 'Item']:
        encoders[col] = joblib.load(f'models/encoder_{col}.joblib')
```

```

if area not in encoders['Area'].classes_:
    print(f'Error: '{area}' not found in training data.')
    return None, 0 # Return early to prevent undefined variables

if item not in encoders['Item'].classes_:
    print(f'Error: '{item}' not found in training data.')
    return None, 0 # Return early to prevent undefined variables

# Now, encode the categorical variables safely
area_encoded = encoders['Area'].transform([area])[0]
item_encoded = encoders['Item'].transform([item])[0]

# Create a single sample in array format
sample = np.array([[area_encoded, item_encoded, year, rainfall, pesticides,
avg_temp]])

print("Before Scaling:", sample)
sample_scaled = scaler_X.transform(sample) # Keep only one transformation
print("After Scaling:", sample_scaled)

# For Bi-GRU, we need historical data to form a sequence
# Here we'll duplicate the sample to create a sequence of required length
time_steps = 3
sample_seq = np.expand_dims(sample_scaled, axis=0) # Reshape for Bi-GRU

# Get Bi-GRU prediction
# Fix Bi-GRU input shape
time_steps = 3 # Set time steps as expected by your model
sample_seq = np.tile(sample_scaled, (time_steps, 1)) # Duplicate the sample 3
times
sample_seq = np.reshape(sample_seq, (1, time_steps, 6)) # Reshape for Bi-GRU

# Predict using Bi-GRU
bigru_pred = bigru_model.predict(sample_seq)

# Combine with tabular data for hybrid model
hybrid_features = np.column_stack((bigru_pred, sample_scaled))

```

```

# Get hybrid model prediction
y_pred_hybrid = hybrid_xgb.predict(hybrid_features)
print("Bi-GRU Prediction Shape:", bigru_pred.shape)
print("Hybrid Features Shape:", hybrid_features.shape)

# Calculate confidence score (using a simpler approach)
# Option 1: Using feature importance-based approach
try:
    # Create DMatrix properly
    dmatrix = xgb.DMatrix(np.array(hybrid_features, dtype=np.float32))

    confidence = hybrid_xgb.get_booster().predict(dmatrix, pred_contribs=True)
    confidence_score = min(0.95, max(0.99, np.mean(np.abs(confidence)) / 10))
except Exception as e:
    print(f"Error calculating confidence with XGBoost API: {e}")
    # Fallback method: Use prediction value as a basis for confidence
    prediction_range = [0, 100000] # Approximate range of possible yield values
    normalized_pred = (y_pred_hybrid[0] - prediction_range[0]) /
(prediction_range[1] - prediction_range[0])
    confidence_score = min(0.95,
                           max(0.65, 0.80 + normalized_pred * 0.15)) # Base confidence of
0.80 with adjustment

if year < 1990 or year > 2100: # Adjust the range as per your dataset
    print(f"Warning: Year {year} is out of range. Adjusting to default (2025).")
    year = 2025 # Set a default year or handle it accordingly

# Inverse transform prediction to original scale
prediction = float(scaler_y.inverse_transform(y_pred_hybrid.reshape(-1,
1)).flatten()[0])
confidence_score = float(confidence_score)

return prediction, confidence_score

def get_suggestions(area, item, predicted_yield, rainfall, pesticides, avg_temp):
    """Generate suggestions to improve or maintain crop yield"""
    suggestions = []

```

```

# Basic suggestions based on crop type
crop_suggestions = {
    'Wheat': ['Ensure proper spacing between plants (15-20 cm)', 'Apply nitrogen
fertilizer during tillering stage',
        'Consider drought-resistant varieties for low rainfall areas'],
    'Rice, paddy': ['Maintain proper water level in paddy fields',
        'Consider SRI (System of Rice Intensification) method',
        'Apply organic matter before transplanting'],
    'Maize': ['Plant in rows 75 cm apart with 20-25 cm between plants', 'Apply
fertilizer in split doses',
        'Control fall armyworm with appropriate measures'],
    'Potatoes': ['Ensure proper hilling to prevent greening',
        'Control for late blight disease with appropriate fungicides',
        'Use certified seed potatoes for better yield'],
    'Cassava': ['Plant at the beginning of rainy season', 'Harvest at optimal maturity
(8-12 months)',
        'Control cassava mosaic disease through resistant varieties'],
    'Soybeans': ['Inoculate seeds with Rhizobium bacteria', 'Control for pod-sucking
insects',
        'Maintain adequate soil pH (6.0-6.5)'],
    'Sweet potatoes': ['Use vine cuttings for planting', 'Maintain adequate soil
moisture',
        'Rotate crops to prevent soil-borne diseases'],
    'Sorghum': ['Plant when soil temperature reaches 15°C', 'Control for sorghum
midge and head smut',
        'Consider bird-resistant varieties in affected areas'],
    'Yams': ['Use minisett technique for propagation', 'Provide support for climbing
vines',
        'Control nematodes through crop rotation'],
    'Plantains and others': ['Apply mulch to conserve moisture', 'Ensure proper
drainage to prevent waterlogging',
        'Control black Sigatoka disease with appropriate measures']
}

```

```

# Add crop-specific suggestions (up to 3)

```

```

if item in crop_suggestions:

```

```

    suggestions.extend(crop_suggestions[item][:3])

```

```

# Region-specific suggestions

```

```

tropical_regions = ['India', 'Brazil', 'Indonesia', 'Thailand', 'Malaysia', 'Bangladesh']

```

```

temperate_regions = ['France', 'Germany', 'United Kingdom', 'Canada', 'Ukraine']
arid_regions = ['Egypt', 'Saudi Arabia', 'Algeria', 'Libya', 'Sudan']

if area in tropical_regions:
    suggestions.append('Consider intercropping with nitrogen-fixing legumes to
improve soil fertility')
elif area in temperate_regions:
    suggestions.append('Monitor winter temperatures for potential frost damage to
early plantings')
elif area in arid_regions:
    suggestions.append('Implement drip irrigation to conserve water in this arid
region')

# Rainfall suggestions
if rainfall < 500:
    suggestions.append('Increase irrigation frequency as rainfall is below optimal
level for most crops')
elif rainfall > 1500:
    suggestions.append('Ensure proper drainage to prevent waterlogging and root
diseases')

# Temperature suggestions
if avg_temp < 15:
    suggestions.append('Consider cold-resistant varieties for better yield in cooler
temperatures')
elif avg_temp > 30:
    suggestions.append('Implement shade structures or mulching to reduce heat
stress during peak temperatures')

# Pesticide suggestions
if pesticides < 10000:
    suggestions.append('Monitor pest populations closely as pesticide usage is
relatively low')
elif pesticides > 100000:
    suggestions.append('Consider integrated pest management to reduce chemical
dependence and environmental impact')

return suggestions[:5] # Limit to top 5 suggestions for better UI
# Return top 5 suggestions for better UI

```



```

def create_confidence_graph(confidence):
    """Create an enhanced confidence meter graph"""
    plt.style.use('ggplot')
    fig, ax = plt.subplots(figsize=(8, 2.5))

    # Define color regions
    colors = ['#e74c3c', '#f39c12', '#27ae60']

    # Create a background with color regions
    for i, color in enumerate(colors):
        start = i / 3
        ax.barh(0, 1 / 3, left=start, color=color, alpha=0.3)

    # Create main confidence bar
    ax.barh(0, confidence, color='#2980b9', height=0.5, alpha=0.8)

    # Add confidence value
    ax.text(confidence, 0, f'{confidence:.2f}', va='center', ha='center',
            fontweight='bold', color='white', bbox=dict(boxstyle='round,pad=0.3',
            fc='#2980b9', ec='none'))

    # Add labels for confidence regions
    ax.text(1 / 6, -0.5, 'Low', ha='center', va='center', fontsize=9)
    ax.text(3 / 6, -0.5, 'Medium', ha='center', va='center', fontsize=9)
    ax.text(5 / 6, -0.5, 'High', ha='center', va='center', fontsize=9)

    # Set limits and remove axes
    ax.set_xlim(0, 1)
    ax.set_ylim(-1, 1)
    ax.set_yticks([])
    ax.spines['top'].set_visible(False)
    ax.spines['right'].set_visible(False)
    ax.spines['left'].set_visible(False)

    # Add a title
    ax.set_title('Prediction Confidence Level', fontsize=12, pad=10)

```

```

# Add benchmark indicators
for x in [1 / 3, 2 / 3]:
    ax.axvline(x, color='white', linestyle='-', alpha=0.5, lw=2)

# Convert plot to base64 string for embedding in HTML
buf = BytesIO()
plt.tight_layout()
plt.savefig(buf, format='png', dpi=120)
buf.seek(0)
img_str = base64.b64encode(buf.read()).decode('utf-8')
plt.close()

return img_str

# Routes
@app.route('/')
def index():
    return render_template('index.html')

@app.route('/register', methods=['GET', 'POST'])
def register():
    if request.method == 'POST':
        # Get form data
        name = request.form['name']
        email = request.form['email']
        password = request.form['password']
        confirm_password = request.form['confirm_password']
        role = request.form.get('role', 'farmer')

        # Form validation
        error = None

        if not name or not email or not password:
            error = "All fields are required"
        elif not re.match(r"^[^@]+@[^@]+\.[^@]+$", email):
            error = "Invalid email address"
        elif password != confirm_password:

```

```

        error = "Passwords do not match"
    elif len(password) < 8:
        error = "Password must be at least 8 characters long"

    if error:
        flash(error, 'danger')
        return render_template('register.html')

    # Check if email already exists
    cur = mysql.connection.cursor()
    cur.execute("SELECT * FROM users WHERE email = %s", (email,))
    user = cur.fetchone()

    if user:
        cur.close()
        flash('Email already registered', 'danger')
        return render_template('register.html')

    # Hash password
    hashed_password = generate_password_hash(password)

    # Insert new user
    cur.execute(
        "INSERT INTO users (name, email, password, role) VALUES (%s, %s, %s, %s)",
        (name, email, hashed_password, role)
    )
    mysql.connection.commit()
    cur.close()

    flash('Registration successful, please log in', 'success')
    return redirect(url_for('login'))

return render_template('register.html')

@app.route('/login', methods=['GET', 'POST'])
def login():

```

```

if request.method == 'POST':
    # Get form data
    email = request.form['email']
    password = request.form['password']

    # Validate form
    if not email or not password:
        flash('Please enter email and password', 'danger')
        return render_template('login.html')

    # Check if user exists
    cur = mysql.connection.cursor()
    cur.execute("SELECT * FROM users WHERE email = %s", (email,))
    user = cur.fetchone()

    if user and check_password_hash(user[3], password):
        # Create session
        session['logged_in'] = True
        session['user_id'] = user[0]
        session['name'] = user[1]
        session['email'] = user[2]
        session['role'] = user[4]

        flash(f'Welcome back, {user[1]}!', 'success')
        return redirect(url_for('dashboard'))
    else:
        flash('Invalid email or password', 'danger')

    cur.close()

return render_template('login.html')

@app.route('/dashboard')
def dashboard():
    if 'logged_in' not in session:
        return redirect(url_for('login'))
    return render_template('dashboard.html', areas=AREAS, crops=CROPS)

```

```

@app.route('/predict', methods=['POST'])
def predict():
    # Get form data
    data = request.form
    area = data.get('area')
    item = data.get('item')
    year = int(data.get('year')) # Keep only one conversion

    rainfall = float(data.get('rainfall'))
    pesticides = float(data.get('pesticides'))
    avg_temp = float(data.get('avg_temp'))
    year = int(data.get('year')) # Ensure this is correctly retrieved

    # Validate year range
    if year < 1990 or year > 2100:
        flash("Invalid year selected. Please choose a year between 1990 and 2100.",
            "danger")
        return redirect(url_for('dashboard'))

    # Make prediction
    predicted_yield, confidence = predict_single_sample(area, item, year, rainfall,
        pesticides, avg_temp)

    # Get suggestions
    suggestions = get_suggestions(area, item, predicted_yield, rainfall, pesticides,
        avg_temp)

    # Create confidence graph
    confidence_graph = create_confidence_graph(confidence)

    # Prepare results
    result = {
        'yield': round(predicted_yield, 2),
        'yield_tonnes': round(predicted_yield / 100, 2),
        'tonnes_acres': round(predicted_yield * 0.0000892, 5),
        'confidence': round(confidence, 2),
        'confidence_graph': confidence_graph,
        'suggestions': suggestions,

```

```

        'input_data': {
            'area': area,
            'item': item,
            'year': year,
            'rainfall': rainfall,
            'pesticides': pesticides,
            'avg_temp': avg_temp
        }
    }

    return jsonify(result)

@app.route('/profile')
def profile():
    if 'logged_in' not in session:
        flash('Please log in first', 'danger')
        return redirect(url_for('login'))

    # Get user details
    cur = mysql.connection.cursor()
    cur.execute("SELECT * FROM users WHERE id = %s", (session['user_id'],))
    user = cur.fetchone()
    cur.close()

    if not user:
        flash('User not found', 'danger')
        return redirect(url_for('login'))

    # Create user object to pass to template
    user_data = {
        'id': user[0],
        'name': user[1],
        'email': user[2],
        'role': user[4],
        'created_at': user[5]
    }

    return render_template('profile.html', user=user_data)

```

```
@app.route('/logout')
def logout():
    session.clear()
    flash('You have been logged out', 'info')
    return redirect(url_for('login'))

if __name__ == '__main__':
    app.run(debug=True)
```