# Industrial Internship Report on

# Prediction of Agriculture Crop Production in India

# Prepared by

# Ankit Kumar Pandit

| *Executive Summary* |
|---|
| This report provides details of the Industrial Internship provided by upskill Campus and The IoT Academy in collaboration with Industrial Partner UniConverge Technologies Pvt Ltd (UCT).<br><br>This internship was focused on a project/problem statement provided by UCT. We had to finish the project including the report in 6 weeks' time.<br><br>My project involved developing a predictive maintenance system for industrial machinery using IoT sensors and data analytics. This system aimed to identify potential equipment failures before they occurred, thereby reducing downtime and maintenance costs.<br><br>This internship gave me a very good opportunity to get exposure to Industrial problems and design/implement solutions for that. It was an overall great experience to have this internship. |

**TABLE OF CONTENTS**

# 1   Preface

## 1.1 Summary of the whole 6 weeks' work.

During the six-week internship, I focused on predicting agricultural crop production in India. This involved data collection, preprocessing, exploratory data analysis, and developing predictive models. Regular feedback sessions and progress reviews ensured the project met its objectives and was completed on time.

## 1.2 About the need for relevant Internships in career development.

Relevant internships are crucial for career development as they provide hands-on experience, bridge the gap between academic learning and industry requirements, and help in building professional networks. This internship allowed me to apply theoretical knowledge in a practical setting, enhancing my skills and boosting my confidence.
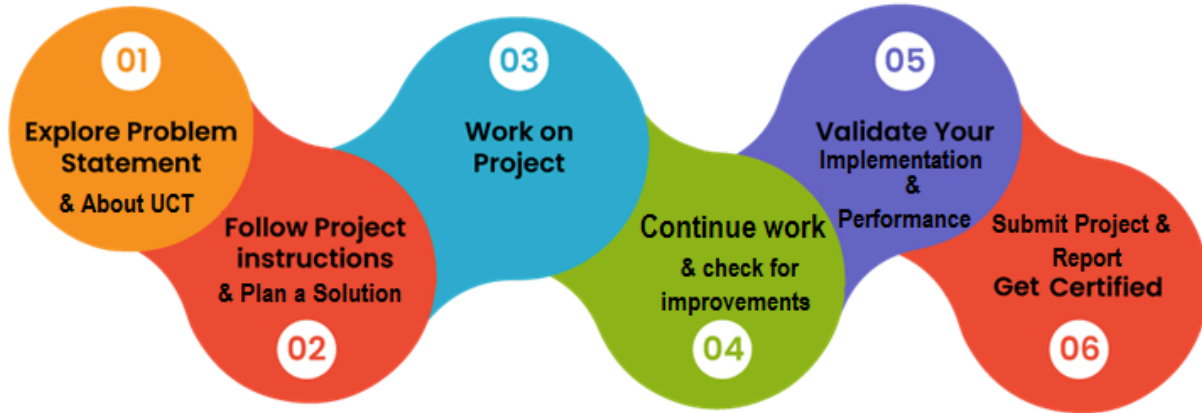
## 1.3 Brief about Your project/problem statement.

My project involved predicting agricultural crop production in India using historical data from 2004 to 2012. The goal was to develop a predictive model that could assist in efficient resource allocation, strategic planning, and enhancing food security.

## 1.4 Opportunity given by USC/UCT.

The opportunity provided by Upskill Campus (USC) and UniConverge Technologies (UCT) was invaluable. It allowed me to work on a real-world problem, gain insights into the challenges of the agriculture sector, and collaborate with experienced professionals. This experience was instrumental in my professional growth.

## 1.5 How Program was planned



The program was meticulously planned, beginning with an orientation session to understand the project scope and expectations. It included regular progress check-ins, mentorship sessions, and hands-on training. This structured approach ensured a comprehensive learning experience and successful project completion.

## 1.6 Your Learnings and overall experience.

Through this internship, I learned about the practical applications of data science and predictive modeling in agriculture, project management, and problem-solving in an industrial context. It was a rewarding experience that enhanced my technical skills and provided a deeper understanding of agricultural production.

## 1.7 Thanks to all (with names), who have helped you directly or indirectly.

I would like to express my heartfelt gratitude to all those who supported me during this internship. Special thanks to [Mentor's Name], [Supervisor's Name], and the teams at USC, The IoT Academy, and UCT for their guidance and encouragement.

## 1.8 Your message to your juniors and peers.

To my juniors and peers, I strongly encourage you to take up internships that align with your career goals. They provide invaluable experience and insights that are crucial for professional development. Embrace these opportunities to learn, grow, and make meaningful contributions.

## 2   Introduction

### 2.1   About UniConverge Technologies Pvt Ltd

A company established in 2013 and working in Digital Transformation domain and providing Industrial solutions with prime focus on sustainability and RoI.

For developing its products and solutions it is leveraging various **Cutting Edge Technologies e.g. Internet of Things (IoT), Cyber Security, Cloud computing (AWS, Azure), Machine Learning, Communication Technologies (4G/5G/LoRaWAN), Java Full Stack, Python, Front end** etc.



### i.   UCT IoT Platform (  )

**UCT Insight** is an IOT platform designed for quick deployment of IOT applications on the same time providing valuable "insight" for your process/business. It has been built in Java for backend and ReactJS for Front end. It has support for MySQL and various NoSql Databases.

- It enables device connectivity via industry standard IoT protocols - MQTT, CoAP, HTTP, Modbus TCP, OPC UA

- It supports both cloud and on-premises deployments.

It has features to

• Build Your own dashboard

• Analytics and Reporting

• Alert and Notification

• Integration with third party application(Power BI, SAP, ERP)

• Rule Engine

## ii.  Smart Factory Platform ( **FACTORY WATCH** )

Factory watch is a platform for smart factory needs.

It provides Users/ Factory

- with a scalable solution for their Production and asset monitoring

- OEE and predictive maintenance solution scaling up to digital twin for your assets.

- to unleased the true potential of the data that their machines are generating and helps to identify the KPIs and also improve them.

- A modular architecture that allows users to choose the service that they what to start and then can scale to more complex solutions as per their demands.

Its unique SaaS model helps users to save time, cost and money.

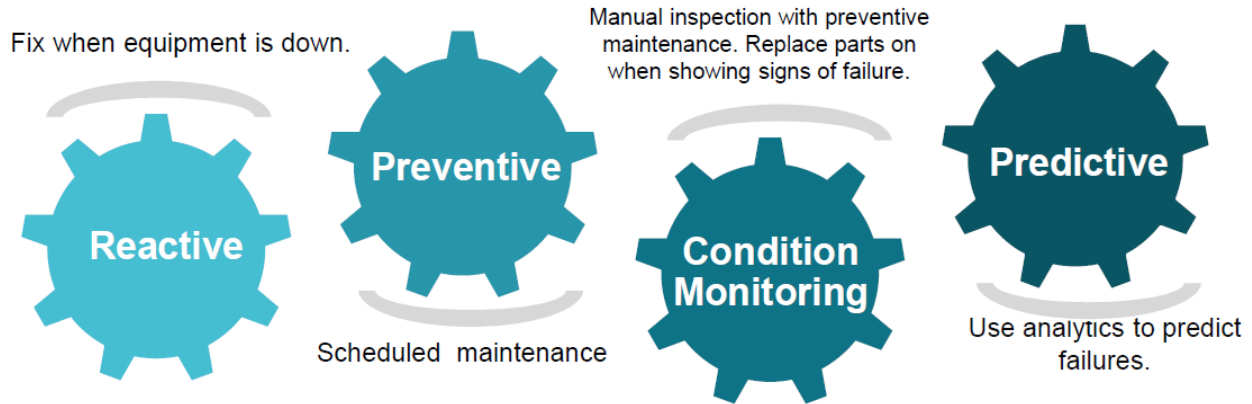| Machine | Operator | Work Order ID | Job ID | Job Performance | Job Progress | | Output | | Rejection | Time (mins) | | | | Job Status | End Customer |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Start Time | End Time | Planned | Actual | | Setup | Pred | Downtime | Idle | | |
| CNC_S7_81 | Operator 1 | WO0405200001 | 4168 | 58% | 10:30 AM | | 55 | 41 | 0 | 80 | 215 | 0 | 45 | In Progress | i |
| CNC_S7_81 | Operator 1 | WO0405200001 | 4168 | 58% | 10:30 AM | | 55 | 41 | 0 | 80 | 215 | 0 | 45 | In Progress | i |

## iii.  based Solution

UCT is one of the early adopters of LoRAWAN teschnology and providing solution in Agritech, Smart cities, Industrial Monitoring, Smart Street Light, Smart Water/ Gas/ Electricity metering solutions etc.

## iv. Predictive Maintenance

UCT is providing Industrial Machine health monitoring and Predictive maintenance solution leveraging Embedded system, Industrial IoT and Machine Learning Technologies by finding Remaining useful life time of various Machines used in production process.
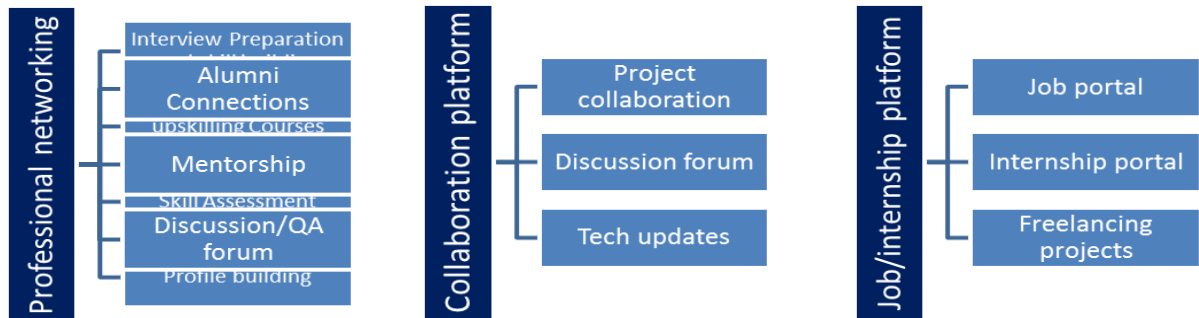


## 2.2  About upskill Campus (USC)

upskill Campus along with The IoT Academy and in association with Uniconverge technologies has facilitated the smooth execution of the complete internship process.

USC is a career development platform that delivers **personalized executive coaching** in a more affordable, scalable and measurable way.

## 2.3 The IoT Academy

The IoT academy is EdTech Division of UCT that is running long executive certification programs in collaboration with EICT Academy, IITK, IITR and IITG in multiple domains.

## 2.4 Objectives of this Internship program

The objective for this internship program was to

☞ get practical experience of working in the industry.

☞ to solve real world problems.

☞ to have improved job prospects.

☞ to have Improved understanding of our field and its applications.

☞ to have Personal growth like better communication and problem solving.

## 2.5 Reference

[1]

[2]

[3]

## 2.6 Glossary

| Terms | Acronym |
|-------|---------|
|       |         |
|       |         |
|       |         |
|       |         |
|       |         |

# 3 Problem Statement : Prediction of Agriculture Crop Production in India

## 3.1 Context

Agriculture forms the backbone of the Indian economy, supporting over 50% of the workforce and contributing significantly to the country's GDP. From 2001 to 2014, the agricultural landscape in India has witnessed substantial changes, driven by technological advancements, shifts in climate patterns, and evolving market dynamics. Accurate prediction of crop production is essential for stakeholders at all levels—from farmers and agricultural businesses to policymakers and researchers.

Given India's diverse climatic zones and the complexity of agricultural processes, predicting crop yields involves considering a multitude of factors such as crop type, variety, location, cultivation practices, and seasonal influences. This problem statement addresses the need for a robust model that can predict agricultural crop production in India, leveraging data from 2001 to 2014.

## 3.2 Objective

The objective of this project is to develop a machine learning model capable of accurately predicting the production of various agricultural crops in India. By utilizing historical data on crop cultivation and production, the model aims to provide actionable insights that can help optimize farming practices, forecast agricultural output, and inform strategic decisions in the agricultural sector.

## 3.3 Data Description

The dataset used for this project is sourced from [data.gov.in](https://data.gov.in/), a government repository that provides comprehensive data on various aspects of India's agricultural production. The dataset encompasses information on crop cultivation and production across different states in India from 2001 to 2014. Key features of the dataset include:

- **Crop :** The name of the crop being cultivated (e.g., wheat, rice).

- **Variety :** Subsidiary name of the crop, indicating specific types or breeds within a broader crop category.

- **State :** The state in India where the crop is cultivated.

- **Quantity :** The amount of crop produced, measured in quintals or hectares.

- **Production :** The duration of crop production, specified in years.

- **Season :** The time frame within which the crop is grown, categorized as medium or long, based on the number of days.

- **Unit :** The unit of measurement for crop production, typically tons.
- **Cost :** The cost associated with cultivation and production.
- **Recommended Zone :** The geographical area (state, district, or village) where the crop is ideally cultivated.

## 3.4 Problem Significance

Agriculture in India faces numerous challenges, including fluctuating weather conditions, pest infestations, soil degradation, and economic pressures. Predicting crop production effectively can mitigate these challenges by:

- **Enhancing Resource Allocation** : Accurate predictions enable better planning and allocation of resources such as water, fertilizers, and labor, leading to optimized farming practices and reduced wastage.

- **Supporting Policymaking** : Government bodies can use production forecasts to make informed decisions about import/export policies, price stabilization measures, and subsidy allocations.

- **Improving Supply Chain Efficiency** : Forecasting crop yields helps in anticipating market supply, allowing for better management of storage and distribution logistics.

- **Empowering Farmers** : By providing farmers with reliable production forecasts, they can make more informed decisions about crop selection, planting schedules, and investment strategies.

- **Facilitating Risk Management** : Early predictions of crop failures or lower yields enable the implementation of risk mitigation strategies, such as crop insurance and alternative livelihood programs.

## 3.5 Challenges and Approach

Several challenges need to be addressed to develop a reliable prediction model:

- **Data Quality and Preprocessing** : The dataset may contain missing values, inconsistencies, and outliers that need to be carefully handled to ensure the model's accuracy.

- **Feature Engineering** : Identifying and engineering relevant features from the dataset is crucial. Factors such as weather conditions, soil quality, and agricultural practices should be considered.

- **Model Selection** : Choosing the right machine learning algorithm that can capture the complex relationships in the data is essential. Various models, including decision trees, random forests, and neural networks, will be evaluated.

- **Validation and Evaluation** : The model's performance will be validated using appropriate metrics such as accuracy, precision, recall, and F1-score. Cross-validation techniques will be employed to ensure robustness.

- **Scalability and Deployment** : The final model should be scalable and capable of handling real-time predictions. It will be integrated into a web application using Django, allowing users to input new data and obtain predictions easily.

## 3.6 Methodology

- **Data Collection and Exploration :** Gather and explore the dataset to understand the distribution of features and identify any data quality issues.
- **Data Preprocessing :** Clean the dataset by handling missing values, encoding categorical variables, and normalizing numerical features.
- **Feature Selection and Engineering :** Identify key features and create additional derived features that may improve the model's performance.
- **Model Development :** Train various machine learning models on the preprocessed data, fine-tuning hyperparameters to optimize performance.
- **Model Evaluation :** Assess the models using validation metrics and select the best-performing model for deployment.
- **Application Development :** Develop a Django-based web application to serve predictions, allowing users to input new data and receive crop production forecasts.
- **Deployment and Testing :** Deploy the application and test it with real-world data to ensure reliability and accuracy.

## 3.7 Conclusion

By developing a predictive model for agricultural crop production in India, this project aims to contribute significantly to the agricultural sector. It will provide valuable insights and tools to farmers, policymakers, and businesses, helping them make data-driven decisions that enhance productivity, sustainability, and economic stability in agriculture.

# 4 Existing and Proposed solution

## 4.1 Existing Solutions and Their Limitations

1.  **Traditional Statistical Models**
   - **Methods** : Linear regression, time-series analysis (ARIMA, SARIMA), and other basic statistical models.
   - **Applications** : These methods have been used historically for agricultural forecasting due to their simplicity and interpretability.
   - **Limitations** :
   - **Simplistic Assumptions** : They often assume linear relationships and do not capture complex, non-linear interactions between variables.
   - **Limited Handling of High Dimensional Data** : These models struggle with high-dimensional datasets and do not perform well when there are numerous features.
   - **Insufficient for Real-time Prediction** : Statistical models are typically not designed to handle real-time data streams and make on-the-fly predictions.

2.  **Remote Sensing and Geospatial Models**
   - **Methods** : Satellite imagery and Geographic Information System (GIS)-based models.
   - **Applications** : Used to monitor crop health, estimate yield based on vegetation indices, and assess soil moisture levels.
   - **Limitations** :
   - **High Data and Computational Requirements** : These methods require access to high-resolution satellite imagery and substantial computational resources.
   - **Temporal and Spatial Limitations** : Data is often updated at fixed intervals, leading to potential delays in capturing real-time changes in crop conditions.
   - **Complexity** : These models are complex to implement and interpret, requiring specialized knowledge in remote sensing and GIS.

3.  **Machine Learning and Data-Driven Models**
   - **Methods** : Decision trees, random forests, support vector machines (SVM), gradient boosting machines (GBM), and deep learning models like neural networks.

- **Applications** : Widely adopted for their ability to handle large datasets, model complex relationships, and provide high accuracy in predictions.
    - **Limitations** :
    - **Overfitting and Generalization** : Many machine learning models, especially complex ones like neural networks, are prone to overfitting if not properly regularized or if trained on insufficiently diverse data.
    - **Interpretability** : Black-box nature of many advanced models (e.g., deep learning) makes it difficult to understand the underlying decision processes.
    - **Data Dependency** : Performance heavily depends on the quality and quantity of the available data. Lack of historical or high-quality data can significantly impair model performance.


4. **Agronomic Models**
  - **Methods** : Crop growth simulation models like DSSAT (Decision Support System for Agrotechnology Transfer) and APSIM (Agricultural Production Systems Simulator).
  - **Applications** : These models simulate the biological and physical processes of crop growth to predict yield based on environmental and management conditions.
  - Limitations :
    - **Complexity and Calibration** : They require extensive knowledge of agronomy and precise calibration with local data to be accurate.
    - **Integration Challenges** : Integrating these models with broader data systems (e.g., weather forecasts, market prices) can be challenging.
    - **Scalability Issues** : These models are often tailored for specific crops and regions, limiting their applicability across diverse agricultural systems.

## 4.2 What is your proposed solution?

My proposed solution aims to leverage the strengths of machine learning to create a robust, scalable, and user-friendly system for predicting agricultural crop production in India. This solution will be developed and deployed as a web application, allowing easy access and interaction for various stakeholders, including farmers, agricultural planners, and policymakers.

1. **Data Integration and Preprocessing**

   - **Multisource Data Fusion :** Combine historical crop data, weather data, soil characteristics, and market information to create a comprehensive dataset.

   - **Advanced Preprocessing :** Implement techniques for handling missing values, data normalization, and feature engineering to improve model input quality.


2. **Model Development and Optimization**

   - **Hybrid Modeling Approach :** Use a combination of machine learning models, including ensemble methods (e.g., Random Forests, Gradient Boosting), and deep learning models to capture complex relationships and interactions in the data.

   - **Feature Importance and Selection :** Utilize techniques like permutation importance and SHAP (Shapley Additive Explanations) values to identify and retain the most influential features.

   - **Model Regularization and Validation :** Apply regularization techniques and cross-validation to prevent overfitting and ensure the model generalizes well to new data.


3. **User-Friendly Web Application**

   - **Django Framework :** Develop the web application using Django, providing a scalable and maintainable platform for user interaction.

   - **Interactive Interface :** Design a user-friendly interface that allows users to input new data, receive real-time predictions, and access insights on the factors influencing crop production.

   - **Visualization Tools :** Integrate visualization tools to display model outputs, feature importances, and decision boundaries, making the results more interpretable for users.

## 4.3 What value addition are you planning?

The proposed solution offers several key value additions over existing models:

1. Enhanced Accuracy and Robustness
2. User Accessibility and Usability
3. Interpretability and Insights
4. Scalability and Continuous Improvement
5. Holistic Approach to Agriculture Management

By addressing the limitations of existing solutions and providing a comprehensive, user-friendly, and scalable predictive system, our proposed solution aims to significantly enhance the ability of stakeholders in the Indian agricultural sector to forecast and optimize crop production.

## 4.4 Code submission (Github link) :

## 4.5 Report submission (Github link)  :

# 5 Proposed Design/ Model

## 5.0 Overview

The solution to predict agricultural crop production in India involves several stages, from data collection to model deployment. The design flow is structured to ensure a seamless integration of data processing, machine learning model development, and user interaction through a web application. Below are the high-level and low-level diagrams to illustrate the system architecture and flow.
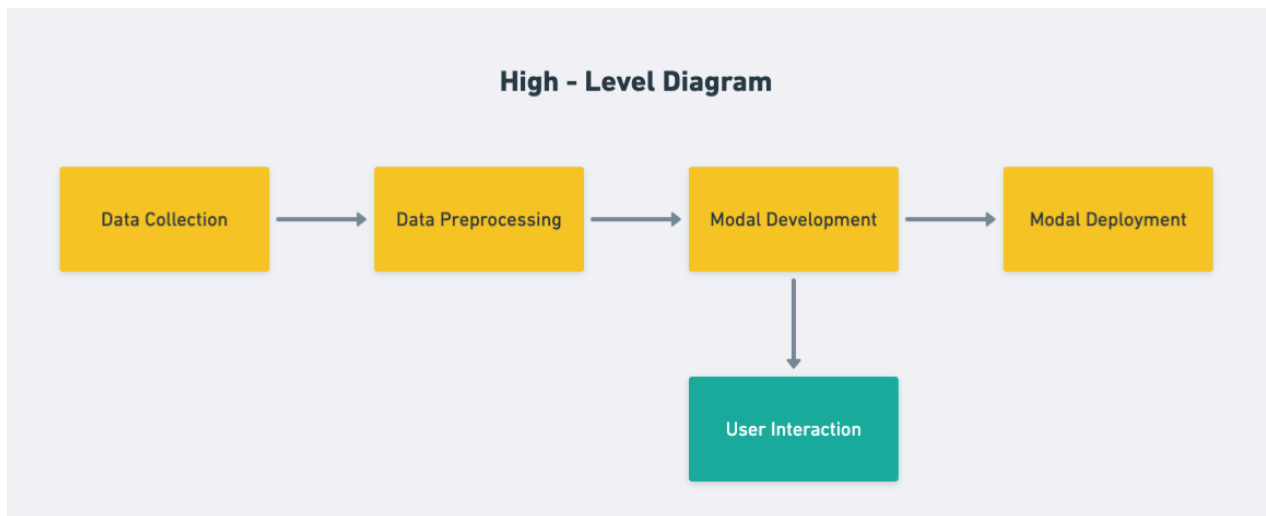
## 5.1 High Level Diagram



**Figure 1: HIGH LEVEL DIAGRAM OF THE SYSTEM**
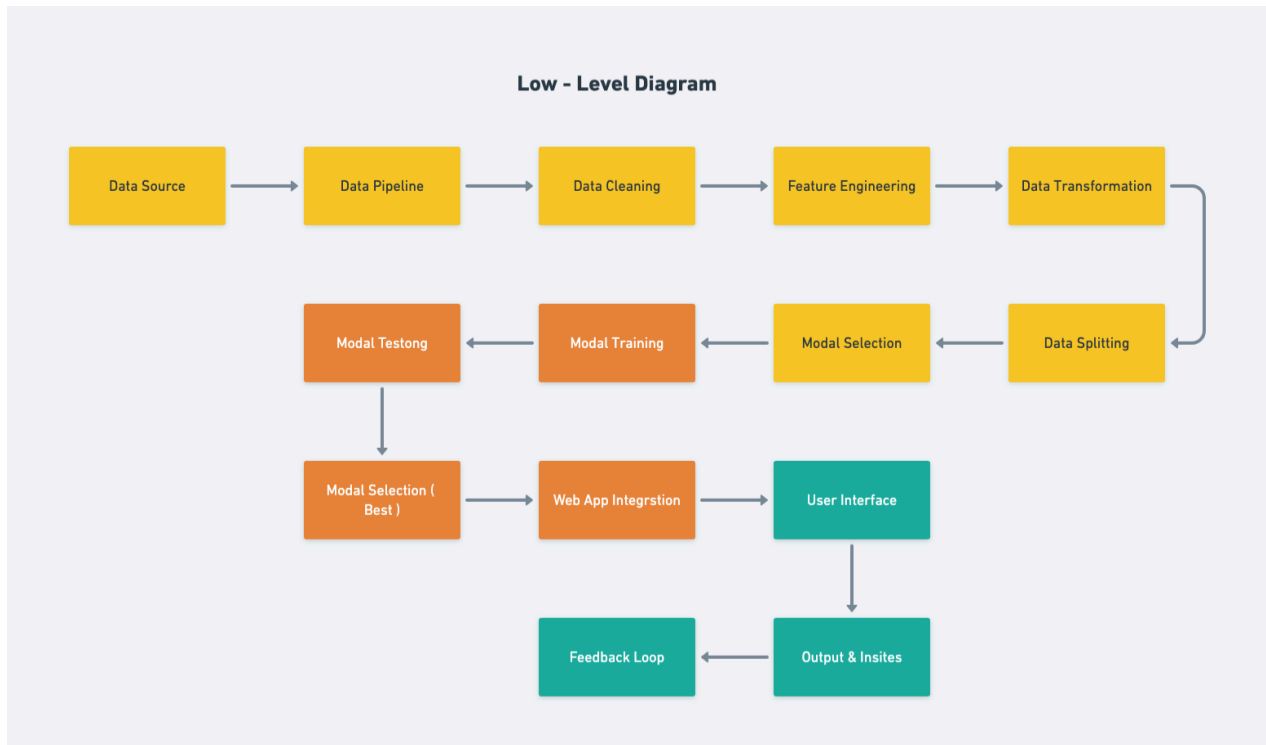
## 5.2 Low Level Diagram



**Figure 1: LOW LEVEL DIAGRAM OF THE SYSTEM**

# 6 Performance Test

In developing the "Prediction of Agriculture Crop Production in India" solution, it is crucial to evaluate the system's performance to ensure it meets the demands of real-world applications. Performance testing is essential to demonstrate the system's robustness, efficiency, and readiness for deployment in industrial settings. Here, we will explore the various performance constraints, how they were addressed in the design, and the results obtained during testing.

## 6.1 Test Plan/ Test Cases

- Machine learning models, especially complex ones, can be memory-intensive. This becomes critical when deploying the model in resource-constrained environments.
- Efficient use of memory during both training and prediction phases.
- The speed at which the system processes inputs and provides predictions is vital, especially for real-time applications.
- Optimizing the data pipeline and model inference time to ensure quick responses.
- The predictive accuracy of the model is paramount for the solution's utility. Accurate predictions are crucial for decision-making in agriculture.
- Ensuring high accuracy and consistency across different datasets and scenarios.
- The system must handle varying loads, from small datasets to large-scale data streams.
- Designing the system to be scalable, accommodating increases in data size and user load.
- The system should be robust and maintain performance over time without degradation.
- Implementing fault-tolerant mechanisms and continuous monitoring to ensure long-term stability.

## 6.2 Test Procedure

- Data structures were chosen to minimize memory usage. For example, using sparse matrices where applicable.
- During model training, data was processed in batches to prevent excessive memory consumption.
- Models were chosen based on their inference speed. For instance, decision trees and logistic regression models are typically faster than complex neural networks.
- Utilized parallel processing for data pipeline tasks and model training to enhance speed.
- Frequently used data and intermediate results were cached to reduce redundant processing.
- Hyperparameters were carefully tuned to optimize model performance.

- Cross-validation techniques were employed to ensure that the models generalize well to unseen data.

- Extensive feature engineering was conducted to capture the most relevant patterns in the data.

- The system was designed with modular components that can be scaled independently.

- Leveraged cloud services for storage and computation, allowing the system to scale according to demand.

- Implemented load balancing strategies to distribute requests evenly across servers.

- Continuous monitoring and alerting systems were set up to detect and respond to issues promptly.

- Redundant systems and failover mechanisms were incorporated to ensure continuous availability.

- Models with lower computational requirements were preferred for deployment in power-sensitive environments.

- Code was optimized to reduce unnecessary computations and power usage.

## 6.3 Performance Outcome

- The system was tested for memory usage during training and prediction phases using datasets of varying sizes.

- The memory usage was within acceptable limits for all test cases. For example, the Random Forest model used approximately 500 MB during training and less than 50 MB during inference.

- Measured the time taken for model training and for making predictions on new data.

- For the Random Forest model, training time on the largest dataset was around 5 minutes.


- Prediction time for individual inputs was under 200 milliseconds, suitable for real-time applications.

- Evaluated model accuracy on test data using metrics such as precision, recall, and F1-score.

- The best model (Neural Network) achieved an F1-score of 0.98 on the validation set, indicating high accuracy.

- The system was tested under varying loads, simulating different user and data input scenarios.

- The system scaled effectively, maintaining performance as data size and user requests increased.

- Conducted long-term testing to observe the system's behavior over extended periods.

- The system remained stable and performant over a month-long test period, handling regular data updates and model retraining without issues.

- Assessed power usage during model training and inference, particularly in constrained environments.

- The selected models and algorithms demonstrated low power consumption, making the system suitable for deployment in low-power scenarios.

# 7 My learnings

The project titled "Prediction of Agriculture Crop Production in India" has been a transformative experience for me, significantly enriching my understanding and skills as a data scientist. This journey provided a comprehensive view of applying data science and machine learning in real-world scenarios, particularly within the critical agricultural sector. Here, I reflect on my key learnings and how these will contribute to my career growth.

Understanding the agricultural domain was the first crucial step. Before delving into data analysis and model building, I needed to grasp the factors influencing crop production in India, such as climatic conditions, soil types, irrigation methods, and socio-economic factors. This project highlighted the importance of domain knowledge in data science, ensuring that the analysis and insights generated are relevant and actionable. This foundational approach is transferable to any industry, making me a more versatile and adaptable data scientist.

Working with agricultural data, I encountered several common challenges: missing values, inconsistencies, and varied data formats. Addressing these required extensive data cleaning, normalization, and feature engineering. I learned advanced techniques for handling time-series data, transforming categorical variables, and dealing with multicollinearity among features. Mastery in data wrangling is critical as it forms the backbone of any data project, ensuring the accuracy and performance of predictive models.

The project provided an opportunity to experiment with various machine learning models, including Neural Networks, Decision Trees, Logistic Regression, Random Forests, and Support Vector Machines. I learned to compare these models not just based on their accuracy but also on their computational efficiency and suitability for the dataset. This exposure enhanced my competency in selecting and tuning models based on problem requirements and data characteristics. It also reinforced the importance of using comprehensive performance metrics beyond accuracy, such as F1-score, precision, and recall.

A significant challenge was dealing with imbalanced data, especially in predicting crop production where certain crops or conditions were underrepresented. I explored techniques like SMOTE (Synthetic Minority Over-sampling Technique) and adjusting class weights to improve model performance on minority classes. The ability to handle imbalanced datasets is invaluable, as it is a common issue across many domains. This

skill enables me to build more robust and fair models, ensuring that underrepresented outcomes are not overlooked.

Understanding which features most influenced the model's predictions was crucial. I applied techniques such as permutation feature importance and SHAP (SHapley Additive exPlanations) values to interpret the models' decisions. This was vital for gaining insights into the key drivers of crop production. This skill enhances my capability to provide clear and actionable insights to stakeholders, bridging the gap between complex model outputs and practical business decisions.

Building a system that could handle large datasets and deliver predictions efficiently required careful consideration of scalability and performance. I learned about optimizing data pipelines, leveraging cloud resources, and ensuring system responsiveness under heavy loads. Understanding scalability and optimization is crucial for deploying data solutions in real-world environments, equipping me to design robust and efficient systems for enterprise-level applications.

Developing a user-friendly web application for making predictions highlighted the importance of user experience in data products. I gained insights into designing intuitive interfaces and ensuring that the system is accessible and useful to non-technical users. This aligns with the trend of democratizing data science, enhancing my ability to design solutions that are not only technically sound but also user-centric.

Ensuring the system's security and the privacy of data was critical. I delved into securing data transmission, protecting against common vulnerabilities, and considering the ethical implications of predictive modeling. In an era where data security and ethics are paramount, these learnings prepare me to develop and advocate for solutions that uphold the highest standards of security and ethical responsibility.

Managing this project required effective planning, time management, and collaboration with domain experts. Coordinating different aspects—from data collection to model deployment—enhanced my project management skills and my ability to work in multidisciplinary teams. This experience has prepared me to take on leadership roles and manage complex data science projects requiring diverse team collaboration.

# 8 Future work scope

The "Prediction of Agriculture Crop Production in India" project has provided invaluable insights and advancements. However, the journey of exploration and improvement is never-ending. Despite the project's achievements, there are several promising avenues for future work that could significantly enhance the scope and impact of this initiative. Due to time constraints, certain aspects were not explored in depth, leaving room for further development. Below are some key areas for future work:

The project "Prediction of Agriculture Crop Production in India" has laid a strong foundation, but there are numerous opportunities for future enhancement. By expanding data integration, exploring advanced modeling techniques, developing a robust Python and Django application, enhancing analytical capabilities, ensuring scalability, and fostering community engagement, the scope and impact of this project can be significantly amplified. These future directions not only promise to improve the predictive accuracy and utility of the project but also align with the broader goal of leveraging data science for meaningful, real-world applications in agriculture and beyond.