# PROJECT REPORT

## ON

## Crop Yield Prediction

A project I Submitted in partial fulfillment of the requirements

for the degree of

B. Tech in CSE (Artificial Intelligence and Machine Learning)

BY

**Aashirwad Wardhan**

**Ayush Taparia**

**Aman Singh**

**Pritam Chowdhury**

**Oishik Mukhopadhyay**

**UNDER THE GUIDANCE OF**

**Dr. Ritamshirsa Choudhuri**

Designation

CSE (Artificial Intelligence and Machine Learning) Department



**TECHNO MAIN SALT LAKE**

EM 4/1 SALT LAKE CITY, SECTOR V

KOLKATA – 700091

West Bengal, India

# Techno Main Salt Lake

[Affiliated by Maulana Abul Kalam Azad University of Technology (Formerly known as WBUT)]

## FACULTY OF CSE (AI & ML) DEPARTMENT

## <u>Certificate of Recommendation</u>

This is to certify that they have completed their project report on **Crop Yield Prediction**, under the direct supervision and guidance of **Dr. Ritamshirsa Choudhuri**. We are satisfied with their work, which is being presented for the partial fulfillment of the degree of B. Tech in CSE (Artificial Intelligence and Machine Learning), Maulana Abul Kalam Azad University of Technology) (Formerly known as WBUT), Kolkata – 700064.

------------------------------------
**Dr. Ritamshirsa Choudhuri**
**(Signature of Project Supervisors)**
**Date:**

------------------------------------
**Dr. Sudipta Chakrabarty**
**(Signature of Project Coordinator)**
**Date:**

------------------------------------
**Dr. Soumik Das**
**(Signature of the HOD)**
**Date:**

# <u>Acknowledgement</u>

We extend our sincere thanks to Dr. Ritamshirsa Choudhuri, our respected professor, for his invaluable guidance and continuous support throughout our project on Time Series Analysis of Stocks and their Predictive Analysis. His mentorship played a pivotal role in shaping our understanding of the subject matter, enabling us to navigate the complexities of the research.

We are grateful to Dr. Soumik Das, Head of the Department, and Dr. Sudipta Chakrabarty, our project supervisor, for providing us with the opportunity to undertake this project. Their trust and encouragement fueled our exploration into the realm of stock prediction and analysis. Their guidance was instrumental in our learning journey, allowing us to delve into the project's nuances.

In the face of challenges encountered during the project, Dr. Ritamshirsa Choudhuri's unwavering support helped us overcome obstacles and shape our ideas into a cohesive project. We express our gratitude to all those who contributed, directly or indirectly, to the successful completion of our research. This project has been a significant learning experience, and we appreciate the collaborative efforts that made it possible.

Date:

| --------------------------- | --------------------------- | --------------------------- |
| Aashirwad Wardhan | Ayush Taparia | Pritam Chowdhury |
| 13030820057 | 13030820007 | 13030820021 |

| --------------------------- | --------------------------- |
| Oishik Mukhopadhyay | Aman Singh |
| 13030820035 | 13030821066 |

# Abstract

India is an agriculture-based nation employing over 50% of the country's workforce. However, despite being called as the backbone of India's economy, the Industry has faced a lot of instability recently. Hence it is only natural that research of various magnitudes is continuously directed towards this field in order to guarantee higher returns in the near future. Keeping in mind the many traditional approaches already prevalent in the agriculture sector which have faced numerous shortcomings, we would like to propose a competent analysis on the Indian agriculture scenario and how climate change is affecting it. This is aimed at providing help to the concerned authorities and any future research. Regression algorithms like Ridge Regression, Random Forest etc. have been used to predict the target variable (Production) and have shown great accuracy within the range of 65%-88 % depending on different algorithms. All this data is deployed on an interactive web application using Streamlit. The proposed system will benefit a lot of people as all information related to India Agriculture can be found out at one place and also the prediction of production amount leads to better planning of resources to be used.

# Table of Contents

# Introduction

In the realm of agricultural science and technology, the prediction of crop yields plays a pivotal role in ensuring food security, sustainable farming practices, and informed decision-making for farmers, policymakers, and stakeholders within the agricultural sector. As the global population continues to burgeon, the demand for food production escalates, making it imperative to enhance the efficiency and productivity of crop cultivation. Crop yield prediction, therefore, emerges as a critical component in addressing the challenges posed by unpredictable climatic conditions, resource constraints, and the need for optimized agricultural practices.

The advent of advanced technologies, including remote sensing, machine learning, and data analytics, has revolutionized the traditional methods of crop yield estimation. This final year project seeks to delve into the development and implementation of a robust crop yield prediction model that harnesses the power of cutting-edge technologies to provide accurate and timely forecasts. By amalgamating agronomic knowledge with computational techniques, the project aims to contribute to the advancement of precision agriculture, offering farmers actionable insights to optimize resource utilization, mitigate risks, and enhance overall agricultural productivity.

Throughout the course of this project, we will explore the intricate interplay between various factors influencing crop yield, such as weather patterns, soil quality, and agronomic practices. Leveraging data-driven approaches, the project endeavors to create a predictive model that goes beyond traditional methodologies, offering a more nuanced and dynamic understanding of the complex dynamics influencing crop growth and yield.

The significance of this endeavor extends beyond the confines of academic exploration. The practical implications of an accurate and reliable crop yield prediction model are vast, ranging from aiding farmers in crop planning and resource allocation to informing policymakers in crafting effective agricultural policies. In the face of a changing climate and increasing demands on global food systems, the development of innovative and precise crop yield prediction models becomes indispensable for fostering sustainable agriculture and ensuring food security on a global scale. This project represents a step forward in harnessing technology for the betterment of agriculture, aligning with the broader goal of creating resilient and efficient food production systems to meet the needs of a growing population.

# Scope of the Project

This project aims to develop an advanced crop yield prediction model, integrating data analytics and cutting-edge technologies. It involves collecting diverse datasets on historical crop yields, weather, soil quality, and agronomic practices for a comprehensive understanding of growth factors.

The project includes an in-depth Exploratory Data Analysis (EDA) and the development of feature engineering techniques, focusing on incorporating remote sensing data. Machine learning algorithms will be selected and fine-tuned for optimal crop yield prediction performance.

Rigorous validation and testing procedures will assess the model's accuracy and reliability. A user-friendly interface will be designed for practical use by farmers and policymakers, facilitating easy access and interpretation of predictions.

Exploration of integration with precision agriculture technologies will enable real-time monitoring and adaptive decision-making in the field. The project emphasizes comprehensive documentation, ethical considerations, and proposes avenues for future research to ensure continuous improvement and adaptation to evolving agricultural needs.

In summary, the main purpose of the project is a comprehensive yet generalized analysis which will act as valuable reference material for future and hence would be easy to use and refer to by concerned authorities and users.The project uses various Machine Learning Algorithms and compares them with each other to find the best one suited for the problem statement.This is a low cost and easy implementation technique when compared to other IOT and image processing models coexisting in the market.The project aims to remove the deficiency of good and comprehensive visualization and interpretation of the entire problem statement. Along with that We are proposing the study of a large amount of production data specific to Indian agriculture production scenario.The project is focussed mainly on the deployment of the project to an Interactive Data Driven Web Application Using Streamlite Which can then be accessed by anyone.The amount of production levels predicted by the model of the project will help farmers to decide the amount of pesticides and other agriculture tools that they might want to use to meet the predicted production requirement (Rajeshkumar J and Kowsigan, M., 2011),lavanya et al 2020.This will help in major cost cutting and save them from the extra payment losses endured in the process.

# Concepts and Problem Analysis

## Time Analysis:

## Gantt Chart

| Team Members | July | August | September | October | November |
|---|---|---|---|---|---|
| Aashirwad Wardhan | Planning | Planning | Research | | Report |
| Ayush Taparia | Planning | | Research | Implementation | Implementation |
| Pritam Chowdhury | Planning | Research | Research | Implementation | |
| Oishik Mukhopadhyay | Planning | Planning | | Report | Report |
| Aman Singh | Planning | Research | Research | Report | Report |

➔ **Planning:** We extensively explored various ideas in collaboration with our professor, Dr. Ritamshirsa, and gained insights into the correct procedure for analyzing crop yield data. This included selecting different types of indicators and algorithms specifically designed for predicting trends in agricultural yields.

➔ **Research**: Having explored diverse analysis techniques, we engaged in thorough research on the algorithms, tools, and methodologies employed by actual agricultural analysts in understanding crop yield data. This involved a comprehensive grasp of these approaches, both in terms of their mathematical foundations and practical applications. Subsequently, we carefully selected a set of these strategies for implementation in our study.

➔ **Implementation:** Python served as our foundational language for all analytical tasks. Leveraging several existing libraries, we conducted the analysis and crafted our implementations tailored to specific requirements. Additionally, we developed a visualization tool using Streamlit, which is instrumental in scrutinizing crop yield data. This tool provides valuable insights, guiding us in decisions related to agricultural practices, such as determining the optimal state for cultivating a specific crop in a given year under certain circumstances.

## Team Structure

In addition to our designated roles, our team operates with a collaborative spirit, fostering open communication and shared decision-making. Regular team meetings will serve as forums for idea exchange, progress updates, and issue resolution, ensuring a cohesive and dynamic workflow. Furthermore, recognizing the dynamic nature of data science projects, we have established flexible timelines and contingency plans to adapt to unforeseen challenges. By leveraging the diverse skills of each team member and maintaining a strong communication framework, we are confident in our ability to not only meet but exceed the objectives outlined in our crop yield prediction project.

In summary, our team, consisting of Aashirwad Wardhan, Ayush Taparia, Pritam Chowdhury, Oishik Mukhopadhyay, and Aman Singh, collaborates seamlessly to tackle the intricacies of crop yield prediction. With Aashirwad, Oishik, and Aman leading research and planning, and Ayush and Pritam focusing on code implementation, we ensure a comprehensive approach to the project. Oishik and Aman will consolidate our findings into a detailed report, capturing the essence of our research journey. Through effective communication, adaptability, and a commitment to excellence, our team aspires to make meaningful contributions to the field of agricultural technology.

## Software Configuration Management:

Software Configuration Management (SCM) is a critical aspect of managing software projects effectively. It involves the systematic control of software artifacts, including source code, documentation, configurations, and dependencies. Here is a suggested approach to SCM for the crop yield prediction project:

**1. Version Control:** Utilize a version control system, such as Git, to manage source code and track changes. Maintain a central repository where team members can collaborate, share code, and track version history. Use branching and merging strategies to manage concurrent development and ensure code stability.

**2. Dependency Management:** Maintain a list of software libraries, packages, and dependencies used in the project. Utilize dependency management tools such as Maven, Pip, or npm to manage and resolve dependencies automatically. Document the versions and configurations of dependencies to ensure consistency across development and deployment environments.

**3. Configuration Management:** Maintain configuration files that capture the project's settings and parameters. Use configuration management tools such as Ansible or Chef to automate the provisioning and management of software configurations across different environments. This ensures consistent behavior and reduces configuration-related errors.

**4. Build Automation:** Automate the build process to ensure reproducibility and consistency. Use build automation tools like Jenkins or Travis CI to build, test, and package the software artifacts. Define build scripts or configuration files that specify the necessary steps to compile, test, and deploy the project.

**5. Release Management:** Establish a release management process to manage the deployment and distribution of software releases. Utilize release management tools to package and distribute the software artifacts, document release notes, and manage versioning. Clearly define release milestones and communicate changes and updates to stakeholders.

**6. Documentation and Change Control:** Maintain documentation that captures the project's design, architecture, and installation instructions. Establish a change control process to track and manage modifications to the software artifacts. Utilize issue tracking systems like Jira or GitHub Issues to capture and prioritize feature requests, bug reports, and change requests.

Adhering to SCM practices ensures proper versioning, reproducibility, and traceability of the software artifacts throughout the project's lifecycle. It also helps in maintaining a stable and manageable codebase, facilitating collaboration among team members, and reducing the risk of errors and inconsistencies.

# Quality Assurance Plan:

A well-defined Quality Assurance (QA) plan is crucial for ensuring the accuracy, reliability, and performance of the crop yield prediction project. Here is a suggested QA plan for the project:

**1. Test Strategy:** Define a comprehensive testing strategy that outlines the various types of testing to be performed, such as unit testing, integration testing, system testing, and acceptance testing. Identify the testing tools and frameworks to be used, and establish guidelines for test case creation, execution, and reporting.

**2. Test Environment:** Set up a dedicated test environment that closely resembles the production environment. This environment should include the necessary hardware, software, and data configurations required for testing. Ensure that the test environment is isolated from the production environment to prevent any impact on live systems.

**3. Test Data:** Prepare a representative and diverse set of test data that covers different scenarios, regions, seasons, and crop types. The test data should include both historical and simulated data to validate the accuracy and robustness of the predictive models. Ensure the test data is anonymized and does not contain any sensitive information.

**4. Test Execution:** Develop and execute test cases based on the defined testing strategy. Test the various components of the project, including data preprocessing, model training, evaluation, and prediction. Conduct both functional and non-functional testing to validate the correctness of the models, accuracy of predictions, and performance under varying conditions.

**5. Bug Tracking and Reporting:** Implement a bug tracking system to capture and prioritize issues identified during testing. Assign severity levels to each bug and track their resolution. Generate regular test reports that summarize the test coverage, test results, and identified issues. Communicate test findings to the development team for prompt resolution.

**6. Documentation:** Maintain comprehensive documentation that describes the testing process, test cases, and test results. Document any issues encountered, their resolution, and lessons learned throughout the project. This documentation serves as a reference for future iterations, enhancements, and maintenance of the crop yield prediction system.

**7. User Acceptance Testing:** Involve end users, such as farmers or agricultural experts, in the testing process. Conduct user acceptance testing to gather feedback on the usability and effectiveness of the system. Incorporate user feedback to improve the user interface, interpretability of results, and overall user experience.

By following a well-structured QA plan, the crop yield prediction project can ensure the reliability and accuracy of the predictive models, enhance user satisfaction, and mitigate risks associated with incorrect predictions or system failures.

# Risk Management:

Effective risk management is essential to identify, assess, and mitigate potential risks that may impact the success of the crop yield prediction project. Here is a suggested approach to risk management for the project:

**1. Risk Identification:** Identify potential risks that may arise during different stages of the project, such as data collection, preprocessing, model development, training, and deployment. Risks could include data quality issues, inaccurate predictions, infrastructure failures, or delays in project timelines. Involve the project team, domain experts, and stakeholders to brainstorm and identify risks.

**2. Risk Assessment:** Assess the identified risks based on their likelihood of occurrence and potential impact on the project. Prioritize risks based on their severity and potential consequences. This assessment helps in allocating appropriate resources and attention to high-priority risks.

**3. Risk Mitigation Strategies:** Develop strategies to mitigate and control identified risks. This may include measures such as data validation and cleaning processes, implementing backup and recovery mechanisms, conducting thorough testing and validation of models, and ensuring scalability and robustness of the deployment infrastructure. Assign responsibilities for risk mitigation activities to the appropriate team members.

**4. Risk Monitoring and Contingency Planning:** Continuously monitor and review identified risks throughout the project lifecycle. Regularly assess the effectiveness of risk mitigation strategies and update them as necessary. Develop contingency plans to address unforeseen risks or events that may impact the project. Maintain open communication channels within the team to promptly address emerging risks.

**5. Stakeholder Communication:** Communicate identified risks, mitigation strategies, and progress in addressing risks to project stakeholders. Keep stakeholders informed about any potential impacts on project timelines, deliverables, or expected outcomes. Engage stakeholders in risk management discussions and decision-making processes.

**6. Documentation:** Maintain a risk register or risk log that documents identified risks, their assessment, mitigation strategies, and outcomes. This documentation serves as a reference for future projects and helps in sharing lessons learned and best practices.

By proactively identifying and managing risks, the crop yield prediction project can minimize potential disruptions, enhance project success, and ensure the reliability and usefulness of the predictive models for farmers.

# Literature Survey

India has been an Agriculture dependent country and not only India most countries rely on their production and agriculture strength. Thus, many studies in this field to predict its future scenario have been performed. . Specific crops like Wheat were found to have a negative relationship with temperature but on the other hand had positive relation with atmospheric $pCO_2$ and rainfall. Also, at higher temperatures there was a sharp incline in rate of decrease in median grain yield when compared to lower temperatures ("Qunying Luo & Amrit Kathuria, 2012"). The Agricultural Production System Simulator Wheat mode (APSIM) is good at considering the physiological effects of $pCO_2$. More IOT based applications have also been researched extensively. Raspberry pi and ML have also been used for a better approach to Agriculture 4.0. Use of Support Vector Machines (SVMs) for embedded systems is considered to be the main commanding classification of engines ("Mahendra Swain, Rajesh Singh, Amit Kumar Thakur and Anita Gehlot ,2019"). Also, for better data collection from different sensors Controller Area Network (CAN) has been used. According to some ("Gbadamosi Babatunde, Adeniyi Abidemi Emmanuel, Ogundokun Roseline Oluwaseun, Oladosu Bukola Bunmi, Anyaiwe Ehiedu Precious, 2018") $CO_2$ Emission and Temperature are not the key factors that affect the production of tuber and root crops. Rather rainfall has more role to play in this scenario. Study of climate factors such as sunlight, soil, humidity etc. was also done to determine a predictive model using regression techniques to find the most optimal temperature and rainfall for crop yield. Many algorithms such as SVM, ERT, RF, DL were compared and their results were analyzed to determine the seasonal sensitivity of corn yield ("Kim Nari & Lee, Yang-Won, 2016"). Also, it was observed that DL methods had the highest accuracies. It was concluded, based on the difference between predictions of USDA statistics and the models used which was about 6-8%, that ML and data mining approaches for prediction and forecasting is a viable option. Talking about different factors affecting Crop production, soil moisture was found to play a major role along with reduced precipitation. Late monsoon and dry season were concluded to be the major reason behind the dry season over central India as well as the resulting decrease in soil moisture. MCA analysis showed that warming of the surrounding oceans such as the Indian Ocean and the Atlantic Oceans is linked to central India and thus resulting in decrease of soil moisture.The effect of temperature is discussed earlier and has worse impact on maize and wheat than on rice, which is more heat tolerant ("Frances C Moore, Uris Lantz C Baldos and Thomas Hertel,  2017"). Also, a comparison between process based and empirical studies shows that there are a lot of differences, based on parameters included which are important. As there have been a lot of changes in both the temperature and the rainfall with the number of days decreasing and the maximum and minimum temperature increased ("Bushra Praveen, Pritee Sharma, 2019") it is important to study these parameters well.Keeping in mind the changing conditions, IOT based implementation as seen before with cloud computing is providing a great way to better predict the agriculture scenario.But at the same time these setups are very difficult to set up and expensive. As the profits in agriculture are very thin and irregular IOT based

implementations can only be used when both funds as well as knowledge to operate these is there ("Olakunle Elijah, Student Member, IEEE, Tharek Abdul Rahman, Member, IEEE, Igbafe Orikumhi, Member, IEEE, Chee Yen Leow, Member, IEEE, and MHD Nour Hindia, Member, IEEE, 2018"). Mann-Kendall test is good for exploring various properties of two basic parameters which are temperature and precipitation as discussed above in other papers. Sen slope is also used in this for the same purpose. After this the cross-distance estimation based on the Euclidean distances between the predicted time series and the station data and this was investigated. Finally, it was concluded that how the level of precipitation and temperature have changed and led to a climate change in the area under study ("Abdul Razzaq Ghumman,Ateeq-ur- Rauf, Husnain Haider and Md. Shafiquzamman ,2019"). Some multi-model ensemble applications using various General circulation models were also studied. This was done for only some major crops like wheat, rice, soybean and maize. For better predictions over the global scenario the mosaic method was useful and efficient. The variability estimation for different crops was done and found out to be within 20%-40%. The final predictions were good and spread up to 40% of the global harvest area. ROC score was also looked at and  visualized properly for each of the four crops ("Toshichika Iizumi, Yonghee Shin, Wonsik Kima, Moosup Kim, Jaewon Choi,2018"). Mining Algorithms like PAM, CLARA etc. have also been used to optimize the predictions and to understand the effect of climate change on the Agricultural scenario. The final result was that the clustering quality of DBSCAN was better than other algorithms like PAM, CLARA etc.("Jharna Majumdar, Sneha Naraseeyappa & Shilpa Ankalaki ,2017"),yasoda et al 2020 and ramamoorthy et al 2020. Climate shock prevention has been studied. The benefits of using crops that are less sensitive to climate change in order to lower the loss that might come if more sensitive crops are planted. Rice was found out to be the most sensitive to precipitation when compared to crops grown in the same season. That is the monsoon season ("Kyle Frankel Davis, Ashwini Chhatre, Narasimha D Rao, Deepti Singh and Ruth DeFries,2019").

# Approach, Methods, and Algorithms

## Introduction:

In the intricate tapestry of agriculture, the accurate prediction of crop yields stands as a pivotal challenge, with implications extending far beyond the farmstead. Traditional methods, tethered to historical data and manual assessments, often fall short in grasping the dynamic nuances of modern agricultural systems. Recognizing this gap, machine learning algorithms have emerged as beacons of innovation, offering the potential to revolutionize crop yield prediction through a data-driven lens. This study focuses on three key algorithms—Random Forest, Decision Tree, and Support Vector Machines (SVM)—as integral components of precision agriculture, aiming to enhance the accuracy and efficiency of forecasting methods.

The Random Forest algorithm, renowned for its ensemble learning approach, amalgamates the strengths of multiple decision trees to create a robust and adaptable predictive model. Decision Trees, with their intuitive and interpretable structure, offer insights into the hierarchical factors influencing crop yields. Meanwhile, Support Vector Machines, as powerful classifiers, bring forth their ability to discern both linear and non-linear relationships within complex agricultural datasets.

To fortify the reliability of these predictive models, the study incorporates cross-validation techniques. This systematic validation methodology involves the iterative division of datasets into subsets for training and testing, ensuring the models' capacity to generalize to unseen data and guarding against overfitting. The amalgamation of advanced algorithms and rigorous validation methods holds the promise of not only enhancing crop yield predictions but also steering agricultural decision-making toward a more sustainable, resilient, and food-secure future.

# METHODOLOGY

## Overview

We are focussed on identifying climate change patterns and its effects in the overall agriculture productivity of India.We are proposing the study of a large amount of production data specific to Indian agriculture production and have tried to predict the future possible effects on agriculture productivity as well.

## Algorithms

We have tried to predict the production of crops specific to the Indian subcontinent based on features like average Temperature, average Rain, State Name, Crop Type and Area to be cultivated.These features have been selected as these are easy to understand and find. These values can be inserted by any user with no technical knowledge and get the desired output easily. We have tried to cover all of the crops that are grown in India. This does add a lot of disappearances and outliers which leads to a lower accuracy. This trade off comes with its own benefits and downfalls. So, they have been taken care of in this project and every aspect of it has been analyzed. Since the target or dependent variable is the Production column, which has continuous values this means we need to go for algorithms that can predict continuous values.Thus, we have gone with regression algorithms for predictions. Regression is a popular technique used for many Machine Learning Projects.

## ALGORITHMS USED

### *Linear Regression Algorithm*

Linear regression, a foundational algorithm in the realm of machine learning, finds valuable application in predicting crop yields within the agricultural domain. At its core, linear regression aims to model the relationship between independent variables (features) and a dependent variable (crop yield) by fitting a linear equation to the observed data. In the context of crop yield prediction, this algorithm proves particularly insightful when attempting to discern linear associations between factors such as weather parameters, soil characteristics, and historical yield data. By estimating the coefficients of the linear equation, linear regression provides a quantitative measure of the impact each variable has on the predicted crop yield.

The simplicity and interpretability of linear regression contribute to its practical utility in agriculture. The algorithm assumes that the relationship between the variables is linear, allowing for a straightforward understanding of how changes in individual factors influence crop yields. However, it is crucial to acknowledge the limitations of linear regression, especially when dealing with complex, nonlinear relationships inherent in agricultural systems. Despite these constraints, linear regression serves as a valuable baseline model, providing a foundational understanding of the relationships within the data.

In the application of linear regression to crop yield prediction, data preprocessing plays a critical role. Handling missing values, normalizing data, and selecting relevant features are essential steps to ensure the model's accuracy. Additionally, incorporating temporal aspects, such as seasonality and climate trends, enhances the algorithm's capacity to capture the dynamic nature of crop growth. While linear regression may not encapsulate the full complexity of agricultural systems, its accessibility, interpretability, and versatility make it an integral component of the toolkit for crop yield prediction, often complementing more sophisticated algorithms in comprehensive agricultural analytics.

# *Random Forest Algorithm*

## Introduction:

In the rapidly evolving landscape of modern agriculture, marked by the relentless pursuit of precision, sustainability, and resource optimization, the integration of advanced machine learning algorithms has become an imperative. Among these, the Random Forest algorithm emerges as a powerhouse, offering a sophisticated ensemble learning approach to predict crop yields. As the agricultural sector grapples with the complexities of climate variability, fluctuating market demands, and the imperative for resource-efficient practices, the Random Forest algorithm plays a pivotal role in revolutionizing predictive modeling. This report embarks on an in-depth exploration of the Random Forest algorithm, elucidating its methodology, strengths, and profound implications within the dynamic context of crop yield prediction in modern agriculture.

## Algorithmic Methodology:

At its core, the Random Forest algorithm harnesses the collective wisdom of multiple decision trees within an ensemble. Each decision tree is constructed using a random subset of the training data and features, mitigating the risk of overfitting and enhancing the model's ability to generalize to diverse agricultural conditions. The iterative amalgamation of these trees, each contributing a unique perspective, results in a robust and adaptable predictive model. This ensemble approach enables Random Forest to capture intricate relationships within the data, making it particularly suitable for the multifaceted variables that influence crop yields, including weather patterns, soil conditions, historical crop performance, and a myriad of other influential factors.

## Practical Implications and Interpretability:

The practical implications of Random Forest in agriculture are far-reaching. Its capacity to aggregate predictions from diverse decision trees not only enhances predictive accuracy but also augments the model's stability and reliability. The algorithm's interpretability extends to the identification of feature importance, offering stakeholders valuable insights into the variables driving crop yield predictions. This transparency is vital for farmers, policymakers, and researchers, providing actionable information for strategic planning, risk mitigation, and the adoption of sustainable agricultural practices. Random Forest's adaptability is further underscored by its seamless integration with various data types, including satellite imagery, historical records, and real-time weather data, positioning it as a pivotal tool in the evolving landscape of precision agriculture.

## Challenges and Considerations:

While Random Forest excels in many aspects, challenges such as computational complexity and the potential for overfitting in certain scenarios must be acknowledged. Balancing model complexity with computational efficiency remains an ongoing consideration. Additionally, the interpretability of Random Forest, while robust, may not match the intuitive simplicity of individual decision trees. Nevertheless, these challenges are often outweighed by the algorithm's predictive power and versatility in addressing the intricate dynamics of modern agriculture.

## Future Directions and Research Avenues:

As technological advancements continue to shape the trajectory of agriculture, the future of Random Forest in the field holds promise. Further research avenues may explore optimizations for computational efficiency, hybrid models combining Random Forest with deep learning techniques, and the integration of emerging technologies like blockchain for enhanced traceability. As agriculture embraces the era of data-driven decision-making, Random Forest stands as a beacon, offering a balance between predictive accuracy and interpretability in the quest for more resilient and sustainable crop yield predictions.

## *Decision Tree Algorithm*

## Introduction:

In the dynamic landscape of contemporary agriculture, the imperative for precise crop yield predictions has intensified, prompting a paradigm shift towards the integration of advanced machine learning algorithms. Among these algorithms, the Decision Tree algorithm stands out as a beacon of transparency and interpretability, providing a robust framework for unraveling the intricate relationships embedded within agricultural datasets. As the agricultural sector grapples with the challenges of climate variability, resource optimization, and the need for sustainable practices, the role of predictive modeling, specifically employing Decision Trees, becomes pivotal in fostering resilience and efficiency. This report undertakes a comprehensive exploration of the Decision Tree algorithm, delving into its methodology, strengths, and practical implications within the ever-evolving landscape of crop yield prediction in contemporary agriculture.

## Algorithmic Methodology:

At its essence, the Decision Tree algorithm embodies a hierarchical structure wherein nodes serve as decision points, contingent upon specific features ranging from weather conditions and soil attributes to historical yield data. The iterative construction of the tree unfolds as the algorithm selectively identifies the most informative features to split the dataset, unraveling granular insights into the myriad factors influencing crop yields. While this methodological approach captures intricate relationships within the data, Decision Trees face the inherent challenge of overfitting. Techniques such as pruning and ensemble methods, most notably in the form of Random Forests, are strategically employed to augment the algorithm's robustness and enhance its capacity to extrapolate patterns beyond the confines of the training data.

## Practical Implications and Interpretability:

The inherent interpretability of Decision Trees becomes a formidable asset in the agricultural sector, where stakeholders demand actionable insights to inform strategic decision-making. The visually traceable decision-making process empowers farmers and other stakeholders to discern the most influential variables affecting crop yields. This transparency proves invaluable in the allocation of resources, risk management, and strategic planning to optimize agricultural practices. Furthermore, Decision Trees can be enriched by the incorporation of additional features such as temporal data and satellite imagery, affording the algorithm the capability to capture the temporal and spatial dynamics integral to crop growth. While Decision Trees may not comprehensively encapsulate the entirety of agricultural complexity, their adaptability and interpretability position them as pivotal tools in unraveling intricate patterns within multifaceted agricultural datasets, steering the course towards the realization of precision agriculture.
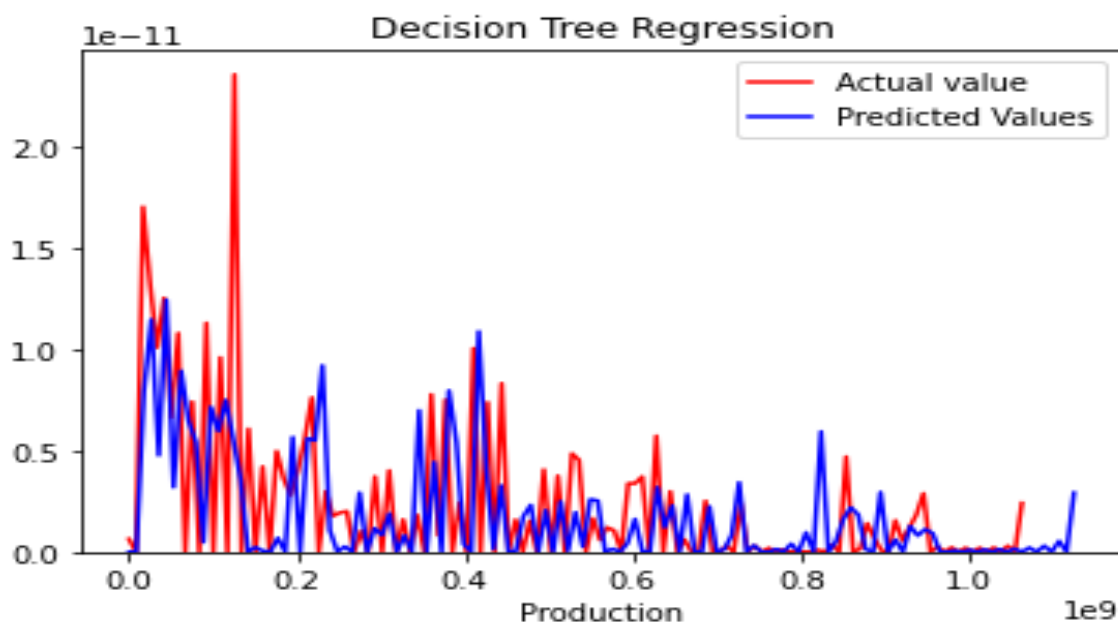
## Challenges and Considerations:

Despite their strengths, Decision Trees are not without challenges. Overfitting, especially in datasets with noise, remains a concern. Pruning helps alleviate this, but finding the optimal balance between complexity and simplicity is an ongoing challenge. Additionally, Decision Trees might struggle with capturing subtle nonlinear relationships, a limitation

addressed by more sophisticated algorithms. However, when integrated into an ensemble, as seen in Random Forests, these challenges are mitigated, making Decision Trees an integral component of comprehensive crop yield prediction models.

## Future Directions:

As technology advances and datasets grow in complexity, the future of Decision Trees in agriculture appears promising. Further research could explore hybrid models, combining Decision Trees with deep learning techniques, to harness the interpretability of the former and the capacity to capture intricate patterns of the latter. Moreover, advancements in remote sensing technology and the integration of real-time data could enhance Decision Trees' ability to adapt to dynamic agricultural conditions. As we navigate the frontier of precision agriculture, Decision Trees continue to be at the forefront, offering a balance between interpretability and predictive power in the quest for more accurate crop yield predictions.

## *Support Vector Regression*

### Introduction:

In the dynamic landscape of precision agriculture, where the integration of advanced machine learning algorithms holds the key to optimizing farming practices, Support Vector Machines (SVM) have emerged as formidable tools for predicting crop yields. This comprehensive report aims to delve into the intricate details of SVM, exploring its algorithmic methodology, strengths, and practical implications in the context of crop yield prediction. As the agricultural sector faces the multifaceted challenges of climate variability, resource optimization, and the imperative for sustainable practices, the role of SVM becomes increasingly pivotal in deciphering the complex relationships within agricultural datasets.

### Algorithmic Methodology:

Support Vector Machines operate by identifying the optimal hyperplane that best separates data into distinct classes. In the context of crop yield prediction, SVM seeks to delineate the intricate relationships between various factors, including weather conditions, soil attributes, and historical yield data. What sets SVM apart is its efficacy in handling both linear and non-linear relationships, making it well-suited for the nuanced and often non-linear interactions within agricultural systems. Additionally, SVM introduces the concept of kernels, enabling the transformation of data into higher-dimensional spaces, thereby enhancing its capacity to discern complex patterns inherent in agricultural datasets.

### Practical Implications and Interpretability:

The practical implications of employing SVM in agriculture are profound, as its capability to handle high-dimensional datasets and discern intricate relationships contributes significantly to predicting crop yields accurately. SVM's adaptability to non-linear relationships proves advantageous in scenarios where traditional linear models may fall short, providing a more nuanced understanding of the variables influencing crop yields. However, it is essential to acknowledge that the interpretability of SVM can be challenging, given that the decision boundary is often represented in a higher-dimensional space. Despite this challenge, SVM offers valuable insights into the complex interplay of variables influencing crop yields, thereby contributing to more informed decision-making in precision agriculture.

### Integration with Agricultural Data:

Support Vector Machines seamlessly integrate with diverse data types, facilitating the incorporation of satellite imagery, real-time weather data, and historical records. This integration enhances the algorithm's capacity to capture the temporal and spatial dynamics of crop growth, providing a holistic view of the myriad factors influencing yields. The adaptability of SVM to various data modalities positions it as a versatile tool for harnessing the wealth of information available in contemporary agriculture, fostering a more comprehensive understanding of the agricultural landscape.

## Challenges, Considerations, and Future Prospects:

Despite its strengths, SVM is not without challenges. The computational complexity of the algorithm, especially in high-dimensional spaces, can pose constraints. Additionally, parameter tuning and model selection require careful consideration to optimize performance. Future research endeavors may focus on refining SVM for efficiency, exploring hybrid models that combine SVM with other machine learning techniques for enhanced predictive power, and adapting the algorithm to accommodate the evolving data sources in precision agriculture. As agriculture continues its digital transformation, the role of SVM in predicting crop yields is poised to expand, contributing to a more resilient, sustainable, and technologically advanced future for precision agriculture.

# *Gradient Boost Algorithm*

## Introduction:

In the dynamic landscape of precision agriculture, the integration of advanced machine learning algorithms plays a pivotal role in shaping data-driven farming practices. Among these algorithms, Gradient Boosting has emerged as a powerful tool for predicting crop yields with unparalleled accuracy. This comprehensive report delves into the intricacies of Gradient Boosting, exploring its algorithmic methodology, strengths, and practical implications in the context of crop yield prediction. As the agricultural sector grapples with the challenges of climate variability, fluctuating market demands, and the imperative for resource-efficient practices, the role of Gradient Boosting becomes increasingly vital in deciphering complex relationships within agricultural datasets.

## Algorithmic Methodology:

Gradient Boosting operates as an ensemble learning technique that builds a predictive model by sequentially combining weak learners, often decision trees. The algorithm minimizes the errors of the preceding model iteration, focusing on the data points that the previous models misclassified. This iterative process results in a robust and accurate predictive model. For crop yield prediction, Gradient Boosting excels in capturing both linear and non-linear relationships among various factors, including weather conditions, soil attributes, and historical yield data. The adaptability of Gradient Boosting to diverse datasets and its capacity to handle complex interactions make it a valuable asset in the agricultural data analytics toolbox.

## Practical Implications and Interpretability:

The practical implications of employing Gradient Boosting in agriculture are profound. Its ability to iteratively refine predictions by learning from past errors makes it exceptionally effective in capturing the nuances of crop yield influencing factors. Gradient Boosting's flexibility allows it to handle missing data and outliers, enhancing its applicability to real-world agricultural datasets. However, similar to other ensemble methods, the interpretability of Gradient Boosting can be challenging due to the combined influence of multiple weak learners. Despite this challenge, the algorithm provides valuable insights into the complex relationships within the data, contributing to informed decision-making in precision agriculture.

## Integration with Agricultural Data:

Gradient Boosting seamlessly integrates with various data types, enabling the incorporation of diverse sources such as satellite imagery, real-time weather data, and historical records. This integration enhances the algorithm's capacity to capture the temporal and spatial dynamics of crop growth, providing a comprehensive understanding of the multifaceted factors influencing yields. The adaptability of Gradient Boosting to different data modalities positions it as a versatile tool for leveraging the wealth of information available in contemporary agriculture.

## Challenges, Considerations, and Future Prospects:

While Gradient Boosting excels in predictive accuracy, challenges such as potential overfitting and model complexity necessitate careful consideration. Parameters such as learning rate and tree depth require tuning to optimize performance. Future research may explore hybrid models combining Gradient Boosting with other advanced techniques, such as neural networks, for enhanced predictive power. As agriculture continues its digital transformation, the role of Gradient Boosting in predicting crop yields is poised to expand, contributing to a more resilient, sustainable, and technologically advanced future for precision agriculture.

## Dataset

Data for the project has been collected from data.gov.in , a site by the government of India consisting of many other official datasets.The main dataset consists ofCrop Production from 1997 till 2015 (table 1.) state wise with every crop included in the dataset. The dataset has around 2-3 Lakh rows.

| Sl.No | Column | Non-Null | Count | Dtype |
|---|---|---|---|---|
| 0 | State_Name | 240691 | non-null | object |
| 1 | District_Name | 240691 | non-null | object |
| 2 | Crop_Year | 240691 | non-null | int64 |
| 3 | Season | 240691 | non-null | object |
| 4 | Crop | 240691 | non-null | object |
| 5 | Area | 240691 | non-null | float64 |
| 6 | Production | 240691 | non-null | float64 |

The environmental factors considered for the project are rainfall and temperature, as through the literature analysis it was evident that these two factors are the most important since they affect the Indian crops most.The data for average rainfall and temperature was taken from the same website. Both the datasets are from 1901 to 2015. All three datasets were merged with respect to the Year column .The final dataset after merging has values till 2015.

## Cleaning and Removing Outliers

Dataset after merging has to be cleaned and prepared for the machine learning module. We have to deal with both numerical and categorical features. For the numerical features we have to first find the outliers so that we have good data points and their distribution should not be skewed.Through plotting box plots we can easily visualize the distribution of the data in a particular column and can see the outliers. To actually find the exact outliers we have gone for IQR. It is a pretty straight forward approach wherein to find the upper limit we add 1.5*IQRto 75th percentile and for lower bound we subtract 1.5*IQR from 25th percentile.(table 2)

$$IQR = 75th-25th$$

$$UB = (1.5*IQR) + 75th$$

$$LB = 25th - (1.5*IQR)$$

| Column | Upper-bound | Lower-Bound |
|---|---|---|
| Production | 16578.5 | -9825.5 |
| Area | 10928.5 | -6427.5 |

Anything outside the upper bound and lower bound is an outlier and should be removed.We can again plot the box plot to check if the distribution is even or not.

## Normalizing

After removing the outliers we can still see that the Area and Production Columns have really large values when compared to other numerical columns.So it becomes necessary to normalize these values so that the machine learning model can predict better. We have used a min-max scaler to normalize and scale the values before we go for the ML module. Min-Max scaler shrinks the data between 0 and 1 without changing the distribution of the data.

Mathematically min-max scaler can be seen as:

$$x\_std=(x-x(min))/(x(max)-x(min))$$

$$x\_scaled = x\_std * (max - min) + min$$

|  | Production | Area |
|---|---|---|
| Mean | 0.089 | 0.107 |
| Std | 0.166 | 0.186 |
| Min | 0.000 | 0.000 |
| Max | 1.000 | 1.000 |

## Encoding

Now after we have cleaned and normalized the numerical variables present in the dataset. We have to deal with categorical variables. Nominal variables don't have any order like we have in our dataset State Names and Crop Names. These two columns have to be encoded so that they can be taken as independent features in the Machine Learning Model.We have used exact encoding for both the columns. Effect/Sum encoding is like One hot encoding but instead of keeping some rows fully 0 it replaces these 0's with -1 so that we don't get a Sparse matrix. We can see after encoding both the columns we finally get 157 feature columns.

# Evaluation

**R2 Score:** R2 score is the most common evaluation technique , or coefficient of determination that measures the proportion of the outcome's variation explained by the model, and is the default score function for regression methods.

**k-fold cross validation:** It is a great method for evaluating the performance of a model on a dataset. It is one of the most important methods used for this purpose. It basically divides the dataset into k almost equal parts or folds and then trains the model at hand on k-1 folds. It keeps one-fold aside for testing. It is repeated k times. Each time the MSE is calculated and at the end average MSE is given

# Deployment

The deployment of the entire model has been done using a python library called streamlit. It allows easy data optimization,deployment and statistical analysis with minimal amount of code.It also prevents any requirement of prior knowledge of deploying web service frameworks like Django and Flask. This can be extremely useful in building data dashboards especially when the team consists of mostly non- tech members.Streamlite is easy to use since it uses predefined commands to build an interactive data driven web application.Simple commands like st.write() to implement a wide variety of objects right from simple text to pandas dataframe matplotlib visualizations becomes possible.Our deployed model consists of a landing homepage along with the navigation sidebar that the user can interact with . The homepage consists of the basic information about the project and accustoms the user to the entire problem statement.The dataset page displays the project's dataset along with its description in a tabular form using the pandas data frame function .The graphs section consists of various interactive graphs plotted using the matplotlib libraries Lastly the model page comprises the user input parameters which include the crop name, state name , amount of area , temperature and rainfall for predicting the target production levels in ton. The user is also able to select the algorithm that they want to apply for better comparison and accuracy after clicking on the run button. Given below is a pseudo code for the Streamlit data driven web application. I

# RESULTS AND DISCUSSION

## Exploratory Data Analysis

After performing cleaning and normalizing the datasets we performed Exploratory Data Analysis to get a better understanding of the agricultural scenario of India. Through plotting and analyzing the graph I came to know a lot of useful information and could understand the dataset better.



Figure 1. Crop-wise total production from 1997 - 2015

Figure 2. Comparison between rainfall, ph,humidity,temperature,N,P,K



Figure 3. Comparison of crop production and average rainfall and comparison of crop production and average temperature

Figure 3 highlights the comparison of crop production with the average rainfall over the years. It is quite evident from the graph that with an increase in rainfall the productivity levels have subsequently shot up whereas whenever a drop of rainfall is seen the productivity also falls. This graph clearly highlights the high dependency of the Indian crops on rain. Figure 13 deals with the comparison of average temperature with the total crop production over the years. The crop production levels have been inversely related to the temperature levels recorded with very high or very low temperature causing major dips in the production levels.



Figure 4. Comparison of crop production and ph and comparison of crop production and humidity

Figure 4 illustrates the correlation between crop production and soil pH levels throughout the observed years. The graph vividly depicts a discernible pattern wherein crop productivity experiences a notable surge with an increase in soil pH. Conversely, a decline in crop production is evident when there is a decrease in soil pH. This graphical representation underscores the significant impact of soil pH on crop yields, emphasizing the importance of maintaining optimal pH levels for sustained agricultural productivity.

In a similar vein, Figure 4 delves into the comparison of crop production with humidity variations over the years. The graph underscores a clear relationship between humidity levels and crop productivity, showcasing that higher humidity is associated with enhanced crop yields. Conversely, a decrease in humidity corresponds to a decline in crop production. This graphical analysis accentuates the pronounced influence of humidity on the agricultural landscape, providing valuable insights for farmers and policymakers aiming to optimize crop production in varying environmental conditions.

Figure 5. Production of Crop in each state

Figure 5 presents a comprehensive overview of crop production distribution across different states. The graph reveals distinctive patterns in agricultural output, with certain states emerging as key contributors to the nation's overall production. Notably, states with favorable climatic conditions and robust agricultural practices showcase higher crop yields, contributing significantly to the country's food production. Conversely, regions facing challenges such as water scarcity or adverse weather conditions exhibit fluctuations in crop production. This graphical representation serves as a valuable tool for policymakers and stakeholders, offering insights into regional disparities and highlighting areas where targeted interventions may be necessary to enhance overall agricultural productivity and ensure food security across diverse geographical landscapes.

# Results

After our Exploratory Analysis of the dataset we arrived at various facts about the Indian agricultural scenario as mentioned in observation. We trained 4-5 Machine Learning Models on the final dataset. We have calculated the R2 score and results are as follows

## Performance of the Algorithms

➢ **Naive Bayes Classification:** Accuracy of the Algorithm Naive Bayes Classification was 0.9709090909090909. Given below is the confusion matrix of Naive Bayes Classification algorithm.



Figure 6. Confusion Matrix of Naive Bayes Classification

➢ **Decision Tree Classification:** Accuracy of the Algorithm Decision Tree Classification was 0.9563636363636364. Given below is the confusion matrix of Decision Tree Classification algorithm.



Figure 6. Confusion Matrix of Decision Tree Classification

➢ **Random Forest Classification:** Accuracy of the Algorithm Random Forest Classification was 0.98. Given below is the confusion matrix of Random Forest Classification algorithm.



Figure 6. Confusion Matrix of Random Forest Classification

➢ **KNN Classifier:** Accuracy of the Algorithm KNN Classifier was 0.9745454545454545. Given below is the confusion matrix of KNN Classifier algorithm.



Figure 6. Confusion Matrix of KNN Classifier

➢ **Gradient Boosting:** Accuracy of the Algorithm Gradient Boosting was 0.9818181818181818. Given below is the confusion matrix of Gradient Boosting algorithm.



Figure 6. Confusion Matrix of Gradient Boosting

# Visualization

Analyzing each type of Crop

**1. Rice:**



**Observations Obtained:**

1. Rice yield is maximum in Rabi season.
2. Rice yield is maximum in Chandigarh.
3. Rice yield has been growing a little from the year 2009 to 2014.

## 2. Wheat:









Observations Obtained :

1. Wheat yield is maximum in Rabi season.
2. Wheat yield is maximum in Punjab.
3. Coconut yield is decreasing in the year 2012 to 2015

## 3.  Coconut:



**Observations Obtained :**

1. Andhra Pradesh is the largest producing coconut state.
2. Production per unit area is higher in Mizoram and Sikkim.
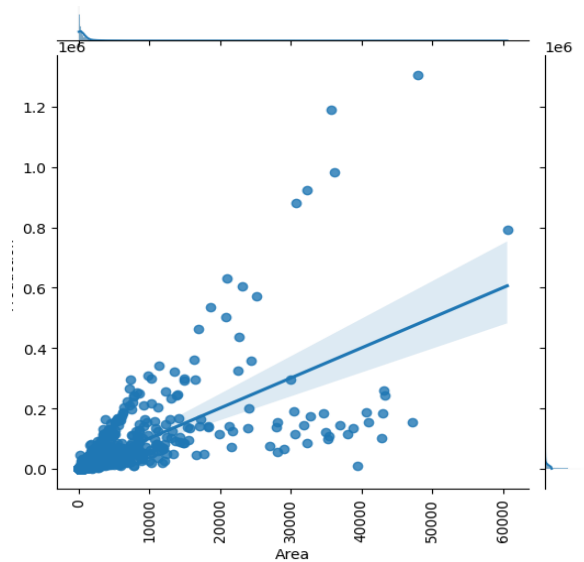3. Coconut yield is decreasing in the year 2012 to 2015.
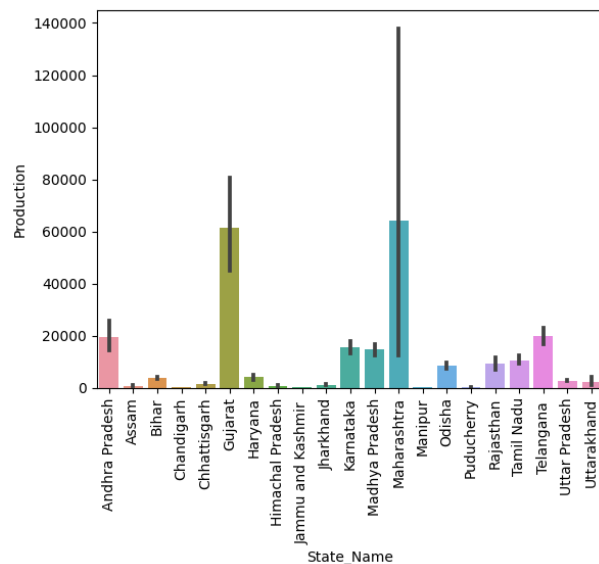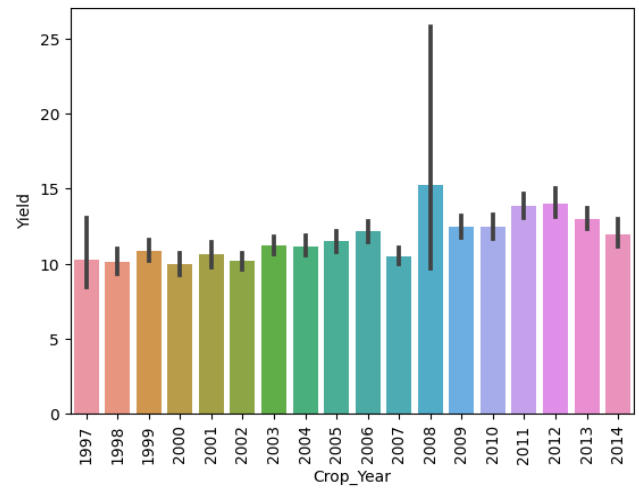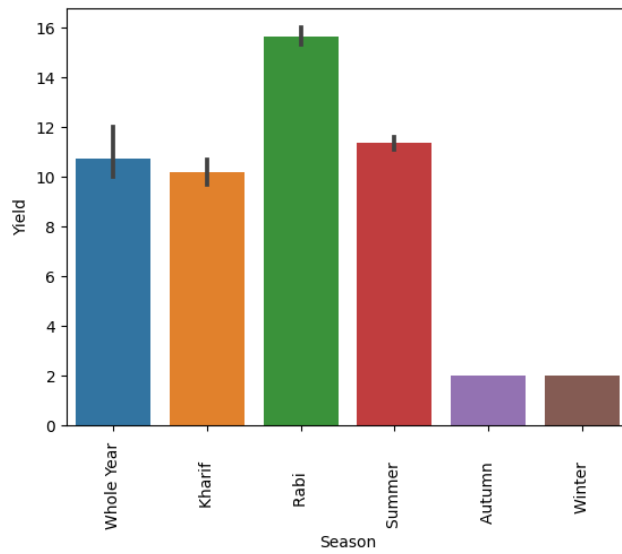
## 4. Potato:









Conclusions Obtained :

1. Potatoes are a Rabi crop.
2. West Bengal is the largest producer of potatoes.
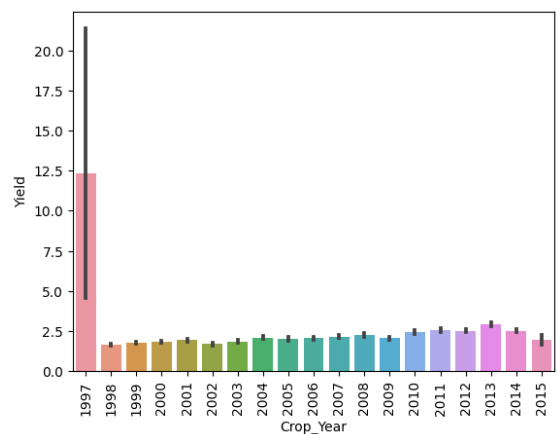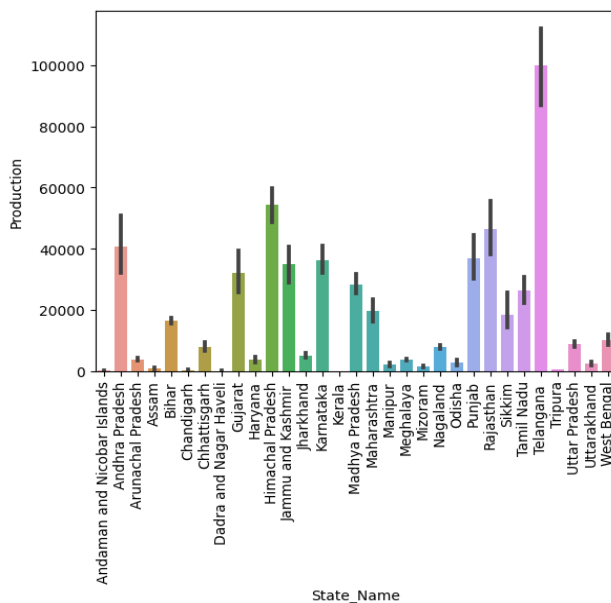3. Potato yield is decreasing in the year 2001 to 2008.
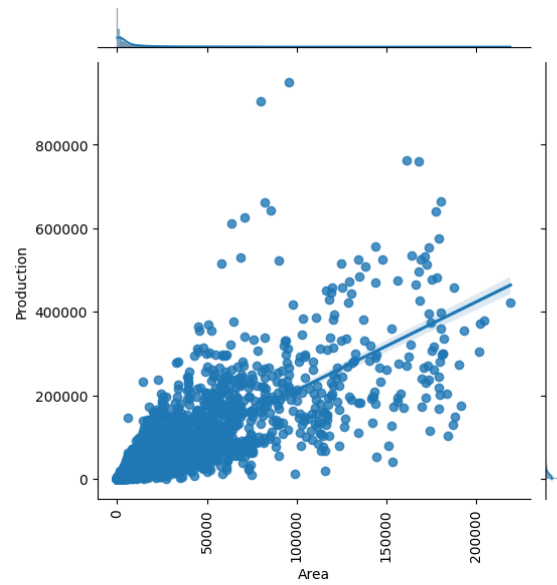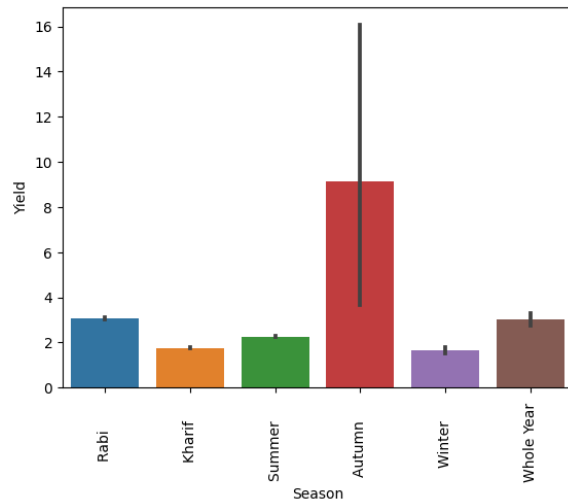
## 5. Onion:



### Observations Obtained:

1. Onion is a Rabi crop.
2. Gujarat and Maharashtra are the major onion-producing states.
3. Onion yield is decreasing in the year 2012 to 2014.

## 6. Maize:



### Observations Obtained:

1. Maize is produced in the autumn season
2. Telangana is the major maize-producing state.
3. There was a sudden decline in maize production from the year 2000.

# Coding Standards and Assumptions

## Coding Standards

**1. PEP 8 Compliance:** The code adheres to the Python Enhancement Proposal 8 (PEP 8) guidelines for code style and formatting. This ensures consistency and readability throughout the codebase.

**2. Modular Structure:** The code is organized into modular components, following best practices for code organization. Each module or function has a clear and well-defined purpose, promoting maintainability and ease of collaboration.

**3. Meaningful Variable Names:** Descriptive and meaningful variable names are used to enhance code readability. This helps in understanding the purpose and role of each variable without additional comments.

**4. Docstrings and Comments:** Comprehensive docstrings are included for functions and modules, providing clear documentation on the purpose, parameters, and return values. Additionally, inline comments are used judiciously to explain complex or critical sections of the code.

**5. Consistent Indentation and Whitespace:** Indentation and whitespace are consistent throughout the codebase, adhering to PEP 8 standards. This ensures a clean and visually coherent code structure.

**6. Version Control Commit Messages:** Clear and informative commit messages are used in version control to track changes effectively. This includes a concise summary of the changes made and, when necessary, additional context or explanations.

## Assumptions

**1. Stationarity of Variables:** The model assumes that the relevant variables influencing crop yield, such as weather patterns and soil conditions, exhibit stationarity over the prediction period. Sudden and unpredictable changes in these variables are not explicitly considered.

**2. Homogeneous Crop Growth:** The model assumes uniform growth conditions across the entire agricultural area under consideration. Localized variations, such as microclimates or specific soil anomalies, are not explicitly incorporated into the prediction.

**3. Consistency in Data Sources:** The assumption is made that data from various sources, such as satellite imagery and historical records, are consistent and accurately represent the corresponding agricultural parameters. Any discrepancies or inaccuracies in the data sources are not explicitly addressed.

**4. Linear Relationships:** In certain models, the assumption is made that the relationships between input variables and crop yield are linear. While this simplification aids in model interpretability, it may not capture more complex, non-linear interactions in the data.

**5. Uniform Crop Management Practices:** The model assumes consistent and uniform crop management practices across the entire agricultural area. Variations in irrigation, fertilization, and other agronomic practices are not explicitly considered.

**6. Predictive Power of Historical Data:** The assumption is that historical yield data provides a reliable indication of future crop yields under similar conditions. Changes in agricultural practices or external factors that might impact yields in novel ways are not explicitly modeled.

It's important to note that these assumptions and coding standards may vary based on the specific requirements of the crop yield prediction model and the nature of the agricultural data being analyzed. Regular validation and refinement of assumptions are recommended as part of the ongoing model development process.

# Conclusion

We can finally conclude the following about the project: Through Exploratory data analysis we could see that Indian crops are directly related to rainfall and negatively to temperature. We could also conclude that the kharif season contributes most to the crop production. The Machine Learning Algorithms Gradient Boost and Random forest Classification had the highest accuracy of about 84% and 84%-76% respectively. Ridge Regression and Gradient Boosting Regression performed fairly well with accuracy of about 54% and 66% respectively. The Data Driven Web application was able to display all the information regarding the project. All the Machine Learning algorithms were successfully deployed on the application. This project signifies the importance of data science and is easy to use in the modern world. The project does come a bit under pressure when compared to specific crop-based applications. It also takes a lot of time to train the complex ML models on this large dataset and thus will require more powerful computers to run it.

The proposed Crop Yield Prediction Algorithm (CYPA) incorporates multiple data sources such as climate, weather, agricultural yield, and chemical data to anticipate annual crop yields by policymakers and farmers in their country. The study demonstrated the efficacy of CYPA by training and verifying five models using optimal hyper-parameter settings for each machine learning technique. The results indicate that CYPA can achieve high accuracy in predicting crop yields, as demonstrated by the scores of Decision Tree Classifier, Random Forest Classifier, and KNN Classifier. Additionally, the paper introduces a new algorithm based on active learning that can enhance CYPA's performance by reducing the number of labeled data needed for training. Incorporating active learning into CYPA can improve the efficiency and accuracy of crop yield prediction, thereby enhancing decision-making at international, regional, and local levels. Overall, this study highlights the potential of IoT and machine learning techniques in addressing the critical challenge of predicting crop yields, thereby facilitating informed decision-making for policymakers and farmers alike. When utilizing Decision Tree Classifier, Random Forest Classifier, and KNN Classifier, the score is equal to 0.95, 0.98, and 0.97, respectively.

# Future Work

This data is ever growing and will require regular updating so that the model is up-to-date. Since this is a huge dataset it would require powerful computers to run more complex algorithms like SVM etc. These algorithms might work better on the dataset and may give better and accurate results. Deep Learning can also be used here in the future to improve the accuracy to a greater extent.

This analysis is just a tip of iceberg, with nineteen year crop production data, a lot could be done and some of the ideas are:

- Instead of deleting missing data for Production(3730 data points), we could impute based on the area used for cultivation and state.
- Zone wise cultivation status and predict future production prediction using regression.
- Crop Categories and status of their cultivation over the years, if the production has gone up (Good case scenario) and if production has gone down (bad case scenario)...can we look into the causation of this trend.
- Asking further important questions like, Kerala is low in area coverage compared to other southern states but still in production levels it's high why?

# References and Bibliography

- https://en.wikipedia.org/wiki/Green_Revolution_in_India

- https://www.kaggle.com/datasets/abhinand05/crop-production-in-india

- Abdul Razzaq Ghumman, Ateeq-ur-Rauf, Husnain Haider and Md. Shafiquzamman , "Functional data analysis of models for predicting temperature and precipitation under climate change scenarios", Journal of Water and Climate Change,2019.

- Bhadouria R, Singh R, Singh VK, Borthakur A, Ahamad A, Kumar G, Singh P (2019) Chapter 1 - agriculture in the era of climate change: consequences and effects. In: Choudhary KK, Kumar A, Singh AK (eds) Climate change and agricultural ecosystems. Woodhead Publishing, Sawston, pp 1–23

- Xu X, Gao P, Zhu X, Guo W, Ding J, Li C, Zhu M, Wu X (2019) Design of an integrated climatic assessment indicator (ICAI) for wheat production: a case study in Jiangsu Province, China. Ecol Indic 101:943–953.

- Alpaydin E (2010) Introduction to machine learning, 2nd edn. MIT Press, Cambridge.