

CROP RECOMMENDATION SYSTEM USING MACHINE LEARNING

Bachelor of Technology
in
Computer Science and Engineering

by

Pathlavath Sudeendra (2021BCS-052)



विश्वजीवनामृतं ज्ञानम्

**ABV INDIAN INSTITUTE OF INFORMATION TECHNOLOGY
AND MANAGEMENT
GWALIOR - 474015**

JULY 2024

CANDIDATES DECLARATION

I hereby certify that the work, which is being presented in the report, entitled **CROP RECOMMENDATION SYSTEM**, in partial fulfillment of the requirement for Summer Colloquium 2024 for **Bachelor of Technology in Computer Science and Engineering** and submitted to the institution is an authentic record of my own work carried out during the period *May 2024* to *July 2024* under the supervision of **Prof. Shashikala Tapaswi**. I also cited the reference about the text(s)/figure(s)/table(s) from where they have been taken.

Date:

Name:

Signature of the Candidate

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Date:

Signature of the Supervisor

ABSTRACT

Agriculture plays a crucial role in India's economy, contributing significantly to the country's GDP and providing employment to a large portion of the population. Despite advancements in agricultural technology, many farmers still rely on traditional methods and make farming decisions based on weather cues. This can lead to suboptimal crop choices and reduced yields.

To address this challenge, this project introduces a machine learning-based Crop Recommendation System that aids farmers in selecting the most suitable crop for their land. The system leverages various machine learning classifiers, including Random Forest, Logistic Regression, Support Vector Machine, and others, to evaluate soil nutrients (like nitrogen, phosphorus, and potassium) and environmental factors (such as temperature, humidity, pH, and rainfall) to recommend the best crop for cultivation.

The system is trained on a dataset containing multiple crops and their respective environmental and soil conditions. The project involves preprocessing steps like feature scaling and standardization, followed by the application of different machine learning models. Among the models tested, the Random Forest Classifier achieved the highest accuracy of 99.3%, making it the primary model for crop prediction. The recommendation system is implemented to provide real-time crop suggestions based on user inputs, potentially offering a valuable tool for farmers and agricultural planners.

This machine learning approach to crop recommendation aims to improve agricultural efficiency by ensuring that crops are well-suited to their growing conditions. By integrating modern technology into farming practices, the system promotes sustainable farming and helps farmers make data-driven decisions, ultimately enhancing productivity and supporting the agricultural sector in India.

ACKNOWLEDGEMENTS

I am grateful to Prof. Shashikala Tapaswi for allowing me to function independently and explore with ideas. I would like to take this opportunity to express my heartfelt gratitude to her not only for her academic guidance but also for her personal interest in the project and constant support as well as confidence-boosting and motivating sessions that proved extremely beneficial and were instrumental in instilling self-assurance and trust. The current work has been nurtured and blossomed mostly as a result of her valuable direction, astute judgment, recommendations, constructive criticism and an eye for perfection. Only because of her tremendous enthusiasm and helpful attitude has the current effort progressed to this point. Finally, I am grateful to the Institution and classmates whose constant encouragement served to renew my spirit, refocus my attention and energy and helped me in carrying out this work.

Pathlavath Sudeendra

TABLE OF CONTENTS

ABSTRACT	2
LIST OF FIGURES	5
1 INTRODUCTION	7
1.1 Context	7
1.2 Objectives	8
2 LITERATURE REVIEW	9
2.1 Background of the Project	9
2.2 System Overview	9
2.3 Related Works	10
3 METHODOLOGY	12
3.1 Data Collection	12
3.2 Data Exploration	13
3.3 Data Preprocessing	13
3.3.1 Data Cleaning	13
3.3.2 Data Balancing	13
3.3.3 Encoding Categorical Variables	14
3.4 Feature Scaling	14
3.4.1 Min-Max Scaling	14
3.4.2 Standardization	14
3.5 Model Exploration	15
3.5.1 Logistic Regression	15
3.5.2 Decision Tree	16
3.5.3 Random Forest	17
3.5.4 Bagging Classifier	17
3.5.5 Gradient Boosting	18
3.5.6 Support Vector Machine (SVM)	19
3.5.7 K-Nearest Neighbors (KNN)	19

TABLE OF CONTENTS

5

4

RESULTS

21

4.1

Experimental Analysis

21

4.2

Performance Comparison

23

4.3

Website Implementation

23

4.3.1

Overview of the Website

23

4.3.2

Technologies Used

23

4.3.3

Website Features

24

5

CONCLUSION

26

5.1

Summary

26

5.2

Future Scope

26

5.3

Limitations

27

5.4

Novelty

28

REFERENCES

28

LIST OF FIGURES

3.1	Dataset	12
3.2	Steps Involved In a Model	15
4.1	Accuracy Comparison	22
4.2	Comparison Table	22
4.3	Model Comparison	23
4.4	Crop recommendation using Certain values	25
4.5	Selection of Model	25

ABBREVIATIONS

ML	Machine Learning
LR	Logistic Regression
DT	Decision Tree
RF	Random Forest
KNN	K-Nearest Neighbor
SVM	Support Vector Machine
PH	Potential of Hydrogen

CHAPTER 1

INTRODUCTION

The Crop Recommendation System represents a significant advancement in agricultural technology, utilizing machine learning (ML) algorithms to enhance crop selection and production efficiency. As agriculture remains a cornerstone of many economies, particularly in countries like India, optimizing crop yields through innovative methods is increasingly crucial. Traditional farming practices, which often rely on historical weather patterns and personal experience, can lead to suboptimal crop choices and reduced productivity. This project addresses these limitations by integrating machine learning techniques to analyze various environmental and soil parameters, offering precise crop recommendations tailored to specific conditions.

Machine learning has revolutionized the field of agriculture by enabling the development of sophisticated models that forecast crop yields with high accuracy. These models use diverse datasets, including soil characteristics, climate patterns, historical yield records, and agronomic practices, to provide comprehensive analyses and predictions. The convergence of machine learning with meteorological data further enhances crop selection by allowing farmers to make informed decisions based on current and forecasted weather conditions. This approach not only improves decision-making but also supports sustainable agricultural practices by optimizing resource use and increasing overall productivity.

1.1 Context

In India, agriculture is a major economic sector and a primary source of livelihood for millions. Despite its significance, the sector often relies on traditional farming methods, which can result in inefficient crop choices and lower yields. Many farmers still depend on personal experience and environmental observations rather than data-driven insights. As global food demands and climate variability increase, there is an urgent need for modern approaches to support more effective agricultural practices.

Machine learning offers a transformative solution by analyzing key factors such as soil nutrients, temperature, humidity, pH levels, and rainfall to provide tailored crop recommendations. This shift towards data-driven agriculture helps farmers optimize their crop selection, manage resources more efficiently, and improve overall productivity. The integration of machine learning with weather forecasting and data mining strategies further enhances crop yield predictions and decision-making processes. By adopting these advanced techniques, the agricultural sector can move towards more sustainable and productive practices, bridging the gap between traditional methods and modern technology.

1.2 Objectives

The primary objective of this project is to create a machine learning-based Crop Recommendation System that can help farmers identify the most appropriate crops to plant based on their specific environmental and soil conditions. The system aims to achieve the following goals:

- **Analyze Key Agricultural Factors:** Use machine learning models to evaluate critical environmental and soil parameters, including nitrogen, phosphorus, potassium levels, temperature, humidity, pH, and rainfall, which are essential in determining crop suitability.
- **Develop an Accurate Predictive Model:** Train and compare various machine learning algorithms to determine the most accurate model for predicting the optimal crop. Special focus will be given to models like the Random Forest Classifier, known for its high accuracy.
- **Implement a User-Friendly Interface:** Develop a web-based application using Flask, allowing users to easily input their data and receive real-time crop recommendations.
- **Support Sustainable Agriculture:** Facilitate the adoption of sustainable farming practices by providing reliable, data-driven crop recommendations, thereby enhancing crop yield and supporting the agricultural sector's overall productivity.
- **Bridge the Gap Between Tradition and Technology:** Make advanced technological tools accessible to farmers, enabling them to integrate modern data-driven insights into their traditional farming practices, ultimately improving decision-making and agricultural outcomes.

CHAPTER 2

LITERATURE REVIEW

2.1 Background of the Project

In earlier versions of crop recommendation systems, there were significant limitations that hindered their effectiveness. Previous projects often overlooked the importance of thorough data preprocessing, leading to noisy and inconsistent data that negatively impacted the accuracy of predictions. Additionally, these earlier systems typically relied on a single machine learning model without exploring or comparing multiple algorithms to determine the most effective approach for crop prediction. As a result, the accuracy of these systems was often suboptimal.

To address these shortcomings, this project implements a more comprehensive data preprocessing pipeline, ensuring that the input data is clean, normalized, and ready for analysis. Furthermore, this project explores and compares multiple machine learning models, including Logistic Regression, SVM, and Random Forest, to identify the best-performing algorithm. The final solution is deployed as a user-friendly web application, offering real-time, accurate crop recommendations based on well-processed data and a carefully selected predictive model.

2.2 System Overview

The crop recommendation system in this project processes agricultural data through several stages. First, the data is preprocessed to clean and normalize inputs like soil nutrients, temperature, humidity, pH, and rainfall. Next, key features are extracted and scaled to prepare the data for machine learning models. These models, including Logistic Regression, SVM, and Random Forest, are then trained to predict the most suitable crops. Finally, the system is deployed as a web application, enabling users to input environmental data and receive crop recommendations in real time.

2.3 Related Works

The literature on machine learning-based crop recommendation systems reveals significant advancements and challenges, as well as promising applications in the agricultural sector. Previous works have primarily focused on utilizing basic data preprocessing techniques. In [4], the authors proposed a crop recommendation system using Random Forest as the machine learning algorithm. Basic data preprocessing methods, such as normalization and handling missing values, were applied. The experimental results showed that the Random Forest model achieved an accuracy of 95%. However, the study suggested that the accuracy could be improved with more advanced preprocessing techniques and a comparison with other machine learning models.

Previous studies on crop recommendation systems have often utilized traditional machine learning techniques with basic data preprocessing. In [1], Bondre et al. proposed a crop yield prediction system using the Random Forest algorithm. However, the study was limited by the use of a smaller set of crop types and rudimentary preprocessing methods. As a result, the Random Forest model achieved an accuracy of only 86.35%. This highlighted the potential for improved accuracy with the application of more advanced preprocessing techniques and a more comprehensive dataset.

Suresh et al. [5] primarily focused on detailing crop information and their nutritional values, such as nitrogen (N), phosphorus (P), and potassium (K) levels. Their research involved analyzing these nutritional parameters within a limited dataset to understand their impact on crop yield. Although the study provided valuable insights into crop nutrition, it was constrained by the dataset's size and scope. Additionally, the research included the design and deployment of a system based on these findings, but it did not extensively explore advanced data preprocessing or compare multiple machine learning models techniques.

Reddy et al. [3] concentrated on ensemble machine learning algorithms in their study published in the International Journal of Scientific Research in Science and Technology. Their research focused on comparing various ensemble methods, such as Decision Trees, Naive Bayes, and Random Forests, to determine the most effective algorithms for crop yield optimization. By visualizing the performance of these algorithms, the study provided valuable insights into selecting the best approach for crop recommendation, showcasing the utility of ensemble methods in enhancing agricultural decision-making.

Garanayak et al. [2] utilized modern machine learning techniques in their research. They focused on a subset of five crops and employed Random Forest Regression, SVM Regression, and Polynomial Regression. By splitting the dataset into training and test sets, they evaluated the performance of these models and achieved an accuracy of 94.7% for rice. Their study compared the accuracy of different regression techniques for crop recommendation, providing valuable insights into the effectiveness of contemporary methods in optimizing crop selection.

CHAPTER 3

METHODOLOGY

The methodology section outlines the systematic approach taken to develop the Crop Recommendation system. It encompasses data collection, preprocessing, model development, training, and evaluation.

3.1 Data Collection

Various datasets relevant to crop recommendation were gathered, each containing environmental and soil data linked to different crop types. The selected dataset includes features such as soil pH, temperature, humidity, rainfall, phosphorus (P), potassium (K), and nitrogen (N). It ensures a balanced distribution across 22 crops, including rice, maize, chickpeas, kidney beans, pigeon peas, moth beans, mung bean, black gram, lentil, pomegranate, banana, mango, grapes, watermelon, muskmelon, apple, orange, papaya, coconut, cotton, jute, and coffee. This dataset comprises 2200 records and was essential for developing models that could accurately recommend the best crop based on specific environmental factors.

	N	P	K	temperature	humidity	ph	rainfall
0	90	42	43	20.879744	82.002744	6.502985	202.935536
1	85	58	41	21.770462	80.319644	7.038096	226.655537
2	60	55	44	23.004459	82.320763	7.840207	263.964248
3	74	35	40	26.491096	80.158363	6.980401	242.864034
4	78	42	42	20.130175	81.604873	7.628473	262.717340
...
2195	107	34	32	26.774637	66.413269	6.780064	177.774507
2196	99	15	27	27.417112	56.636362	6.086922	127.924610
2197	118	33	30	24.131797	67.225123	6.362608	173.322839
2198	117	32	34	26.272418	52.127394	6.758793	127.175293
2199	104	18	30	23.603016	60.396475	6.779833	140.937041

2200 rows × 7 columns

Figure 3.1: Dataset

3.2 Data Exploration

Before diving into preprocessing, a comprehensive exploration of the dataset was carried out to gain a deep understanding of its structure and the distribution of various environmental features such as soil pH, temperature, humidity, rainfall, phosphorus, potassium, and nitrogen. The analysis began by inspecting the initial few rows to get an overview of the data, followed by determining the dataset's overall dimensions. Detailed scrutiny of each feature was performed to uncover insights into data types and identify any potential issues, such as missing values or inconsistencies. A statistical summary was then derived to capture key metrics like central tendencies, variability, and the presence of outliers. To further understand the relationships between different features, a correlation analysis was conducted, and the results were visualized through a heatmap. This exploration provided critical insights into the interdependencies among variables, offering valuable guidance for subsequent preprocessing steps.

3.3 Data Preprocessing

3.3.1 Data Cleaning

The initial step in the preprocessing pipeline was focused on ensuring the dataset's cleanliness. This began with a thorough check for missing values, a frequent challenge in real-world data. Fortunately, no missing values were detected, confirming the dataset's completeness. Next, the dataset was examined for any duplicate entries. Upon identifying duplicates, these were promptly removed to preserve the integrity of the dataset, ensuring that each data point was unique and did not introduce any potential bias into the learning process.

3.3.2 Data Balancing

In machine learning, especially in classification tasks, having a balanced dataset is crucial for training models that can generalize well across all classes. The dataset used for the Crop Recommendation System was inherently balanced, meaning that each of the 22 crop types was equally represented. This balance was critical as it prevented the models from becoming biased towards the more frequently occurring classes, ensuring that the system could accurately recommend crops across different environmental conditions.

3.3.3 Encoding Categorical Variables

The crop labels, which were originally in string format, needed to be converted into a numerical format suitable for machine learning algorithms. This was achieved through label encoding. A dictionary mapping (crop-dict) was created to convert the crop names into corresponding numerical values, which were then added to the dataset as a new column crop-num. This transformation was essential because machine learning algorithms require numerical input, and this step ensured that the crop labels were appropriately encoded for the modeling process.

3.4 Feature Scaling

Feature scaling is a vital preprocessing step in the Crop Recommendation System, ensuring that all input features contribute equally to the model's learning process. In this project, two primary scaling techniques were employed: Min-Max Scaling and Standardization. These techniques were applied to the environmental variables, including soil pH, temperature, humidity, rainfall, phosphorus (P), potassium (K), and nitrogen (N).

3.4.1 Min-Max Scaling

The first step in scaling involved normalizing the data using Min-Max Scaling. This technique adjusts the range of the data so that all features lie between 0 and 1. Min-Max Scaling is particularly useful when dealing with features that have different units or scales. For example, temperature and rainfall may naturally have vastly different ranges, which could potentially skew the results of distance-based algorithms like K-Nearest Neighbors (KNN). By applying Min-Max Scaling, all features were brought into the same range, ensuring that no single feature dominated the model's learning process due to its magnitude.

3.4.2 Standardization

Following Min-Max Scaling, the data was further standardized using StandardScaler. Standardization shifts the data distribution to have a mean of 0 and a standard deviation of 1. This process is crucial for algorithms that are sensitive to the distribution of the data, such as Support Vector Machines (SVM) and Logistic Regression. Standardization ensures that each feature contributes equally to the model, especially when the model relies on assumptions about the data's distribution. By applying this technique, the models were better equipped to converge more quickly during training and achieve higher accuracy.

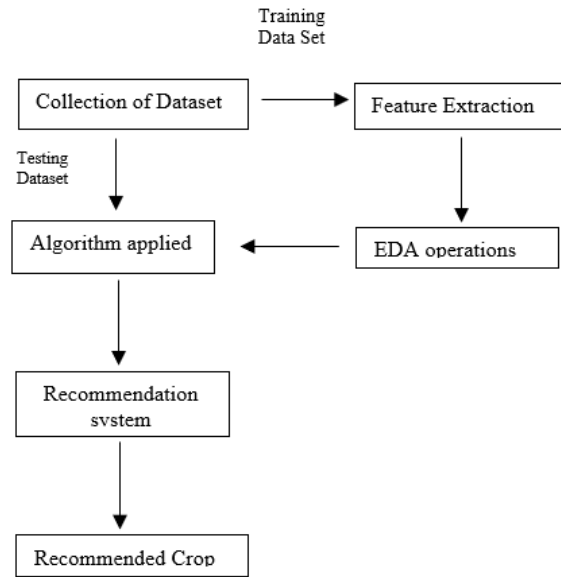


Figure 3.2: Steps Involved In a Model

3.5 Model Exploration

In the Crop Recommendation System, ensemble models like Random Forest and Gradient Boosting are preferred over simpler algorithms due to their ability to capture complex relationships between environmental factors and crop suitability. Factors such as soil pH, temperature, and nutrient levels interact in non-linear ways, which simpler models like Logistic Regression or Decision Trees may not fully capture. Ensemble methods, by combining the predictions of multiple decision trees, provide a more accurate and robust analysis, making them ideal for generating reliable crop recommendations. Random Forest, in particular, stands out due to its high accuracy and effectiveness in handling diverse agricultural data.

3.5.1 Logistic Regression

Logistic Regression is a linear model that estimates the probability of a certain crop being the best recommendation based on environmental features. It's particularly effective when the relationship between features and the target is roughly linear.

1. Advantages

- **Simplicity and Interpretability:** Provides a clear understanding of how each environmental factor, such as temperature or soil pH, impacts the crop recommendation.

- **Quick Benchmarking:** Serves as a baseline model to compare more complex models against.

2. Disadvantages:

- **Linear Relationships:** Assumes a linear relationship between features and the target, which may not capture the intricate interactions between environmental factors and crop yield.
- **Limited Complexity:** Struggles to handle non-linearities and complex patterns in the data, which are common in agricultural datasets.

In the Crop Recommendation System: Logistic Regression provided a starting point to understand the relationships between environmental features and crop recommendations. However, its linear nature limited its ability to capture the complexities of crop growth conditions.

3.5.2 Decision Tree

Decision Tree uses a tree-like structure to make decisions by splitting the data based on feature values. It models decision rules that can directly map environmental conditions to the recommended crop.

1. Advantages

- **Interpretability :** The tree structure provides a clear visual representation of decision rules, making it easy to understand the reasoning behind crop recommendations.
- **No Need for Scaling:** Works well without feature scaling, which simplifies the preprocessing steps.

2. Disadvantages:

- **Prone to Overfitting:** Decision Trees can easily overfit, capturing noise in the data as if it were a true signal, leading to poor generalization.
- **Instability:** Small changes in the data can lead to different trees being formed, making the model less reliable.

In the Crop Recommendation System: Decision Trees provided clear decision paths for crop recommendations, though the risk of overfitting required careful pruning and validation.

3.5.3 Random Forest

Random Forest is an ensemble method that builds multiple decision trees and combines their predictions to improve accuracy and robustness. It aggregates the predictions from each tree to make a final decision on the best crop recommendation.

1. Advantages:

- **High Accuracy:** By averaging the results of multiple decision trees, Random Forest reduces overfitting and enhances predictive accuracy, which is crucial for reliable crop recommendations.
- **Feature Importance:** Provides insights into which environmental features (e.g., nitrogen levels, rainfall) are most influential in determining the best crop, helping to understand the driving factors behind the recommendations.
- **Robustness:** Handles large datasets with higher dimensions effectively and is less sensitive to noisy data and overfitting.

2. Disadvantages:

- **Complexity:** More complex and less interpretable compared to a single decision tree, as it involves multiple models.
- **Computational Resources:** Requires more computational power and memory during training and prediction.

In the Crop Recommendation System: Random Forest was the top-performing model, offering the highest accuracy in predicting the best crop to plant. Its robustness to noise and ability to handle complex interactions between features made it particularly well-suited to this project, ensuring reliable recommendations across diverse environmental conditions.

3.5.4 Bagging Classifier

Bagging Classifier is an ensemble method that builds multiple models from different subsets of the training data and combines their predictions to enhance model stability and accuracy.

1. Advantages:

- **Reduced Variance:** By aggregating the results from multiple models, Bagging reduces the variance of the predictions, leading to more stable and reliable recommendations.

- **Enhanced Accuracy:** Improves overall accuracy by averaging out errors that individual models might make.

2. Disadvantages:

- **Complex Implementation:** More complex to implement compared to individual models, as it requires managing multiple training processes.
- **Increased Computational Load:** Requires training multiple models, which can be computationally expensive.

In the Crop Recommendation System: The Bagging Classifier enhanced the stability and accuracy of crop recommendations by reducing variance, making it a valuable addition to the ensemble methods used in the project.

3.5.5 Gradient Boosting

Gradient Boosting builds models sequentially, with each new model correcting the errors of the previous ones. This method iteratively improves prediction accuracy by focusing on the hardest-to-predict cases.

1. Advantages:

- **High Predictive Power:** Excels at capturing complex patterns in the data by sequentially refining predictions, leading to high accuracy in crop recommendations.
- **Adaptability:** Can be customized to various loss functions and models, making it flexible for different types of prediction tasks.

2. Disadvantages:

- **Risk of Overfitting:** If not properly tuned, Gradient Boosting can overfit the training data, leading to poor generalization.
- **Training Time:** The sequential nature of model building can be slow, particularly with large datasets.

In the Crop Recommendation System: Gradient Boosting provided incremental improvements in accuracy by focusing on correcting errors from previous models. Its ability to handle complex data distributions made it an excellent choice for refining crop recommendations, though careful tuning was required to prevent overfitting.

3.5.6 Support Vector Machine (SVM)

SVM is a powerful algorithm that finds the optimal hyperplane to separate different crop types in a high-dimensional feature space. It is particularly effective in cases where classes are not linearly separable.

1. Advantages:

- **High-Dimensional Capability:** Excels in handling high-dimensional data, making it effective for complex datasets with many features like soil nutrients and weather conditions.
- **Flexibility with Kernels:** The ability to apply different kernel functions allows SVM to model non-linear relationships in the crop recommendation data.

2. Disadvantages:

- **Computationally Intensive:** The training process can be slow, especially with large datasets and when tuning parameters such as the kernel and regularization term.
- **Sensitivity to Parameters:** Requires careful tuning of parameters, and performance can degrade if not properly configured.

In the Crop Recommendation System: SVM was beneficial for its ability to handle the complex and high-dimensional nature of the crop dataset, though its computational demands required careful consideration during model training.

3.5.7 K-Nearest Neighbors (KNN)

KNN is a non-parametric method that classifies a crop by a majority vote of its neighbors. It looks at the closest environmental conditions (e.g., soil pH, temperature) to predict the best crop.

1. Advantages:

- **Non-Linear Decision Boundaries:** KNN can capture complex, non-linear relationships between features and crop types, which is crucial for accurate crop recommendations.
- **Intuitive Understanding:** Simple to understand and implement, as it relies on direct comparisons to similar cases in the training data.

2. Disadvantages:

- **Computational Burden:** Requires significant computation during prediction, especially with large datasets like the crop recommendation system.
- **Sensitive to Noise:** Performance can be negatively impacted by noisy data or irrelevant features, which are common in agricultural datasets.

In the Crop Recommendation System: KNN was effective in capturing the non-linear relationships between environmental factors and crop types, but its computational intensity and sensitivity to noise needed to be carefully managed.

CHAPTER 4

RESULTS

4.1 Experimental Analysis

In our Crop Recommendation System project, we tested several machine learning models to evaluate their performance in predicting the best crops to grow based on environmental factors. The models tested included Logistic Regression, Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Decision Tree, Random Forest, Bagging Classifier, and Gradient Boosting. Each model was assessed based on its accuracy, precision, and recall, providing valuable insights into their effectiveness for this task.

The Logistic Regression model, used as a simple benchmark, achieved an accuracy of 96.34%, with a precision of 96.44% and a recall of 96.36%. Although effective as a baseline, its performance was surpassed by more complex models.

The Support Vector Machine (SVM) showed a slight improvement with an accuracy of 96.82%, precision of 97.15%, and recall of 96.82%. This model's capability to handle complex data was beneficial, but it still did not outperform the ensemble methods.

The K-Nearest Neighbors (KNN) model, known for capturing non-linear relationships, achieved an accuracy of 95.91%, a precision of 96.54%, and a recall of 95.91%. However, its performance depended on the proper selection of parameters.

The Decision Tree model performed well, achieving an accuracy of 98.18%, with a precision of 98.24% and a recall of 98.18%. While it provided clear decision paths, it had a tendency to overfit the data.

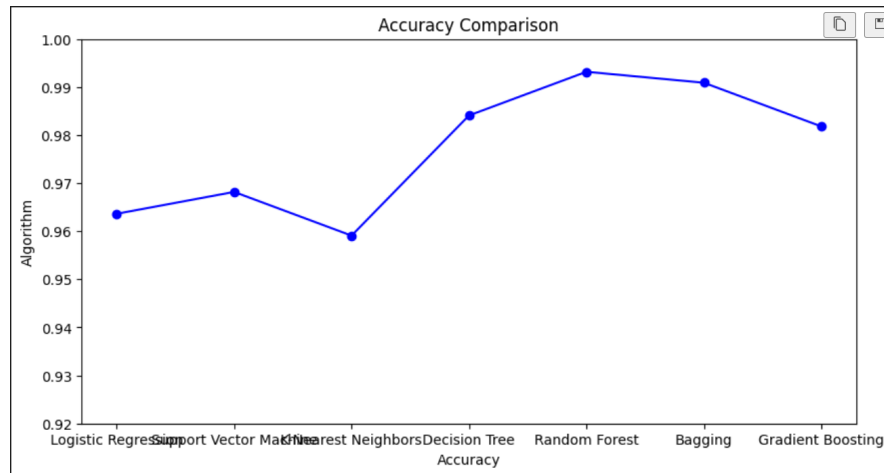


Figure 4.1: Accuracy Comparison

the Gradient Boosting model achieved an accuracy of 98.18%, with a precision of 98.43% and a recall of 98.18%. Its ability to iteratively correct errors of previous models showed promise, though it required more computational resources.

The Bagging Classifier also performed well, with an accuracy of 98.59%, a precision of 97.16%, and a recall of 98.09%. Although it provided enhanced stability, it was slightly less accurate than Random Forest.

Finally, The Random Forest model stood out as the most effective, achieving the highest accuracy of 99.32%, with a precision of 99.37% and a recall of 99.32%. Its strength in handling diverse data and reducing overfitting made it the best choice for this crop recommendation task.

Model	Precision	Recall	Accuracy
Logistic Regression	0.964442	0.963636	0.963401
Support Vector Machine	0.971517	0.968182	0.968182
Bagging	0.971620	0.980909	0.985909
Decision Tree	0.982402	0.981818	0.981818
Random Forest	0.993735	0.993182	0.993182
K-Nearest Neighbors	0.965390	0.959091	0.959091
Gradient Boosting	0.984271	0.981818	0.981818

Figure 4.2: Comparison Table

4.2 Performance Comparison

The performance comparison between the Crop Recommendation System and the previous base model shows a clear improvement in accuracy and reliability, primarily due to the implementation of new techniques. By increasing the number of features, the current model could better capture the complexities of the agricultural environment, leading to more precise crop predictions.

Model	Previous Model Accuracy	Our Model Accuracy
Support Vector Machine	0.75	0.96
K-Nearest Neighbor	0.90	0.95
Random Forest	0.95	0.99

Figure 4.3: Model Comparison

Enhanced data preprocessing techniques, such as careful feature scaling and dataset balancing, played a crucial role in refining the model’s inputs. These improvements helped in reducing noise and ensuring that the model was trained on high-quality data. As a result, the current model demonstrates superior performance, highlighting the effectiveness of these methodological advancements over the earlier approaches.

4.3 Website Implementation

4.3.1 Overview of the Website

The website was developed to provide a user-friendly interface for a crop recommendation system using machine learning. Users can input soil and environmental data such as nitrogen, phosphorus, potassium levels, temperature, humidity, pH, and rainfall, which are then processed by trained models to recommend the most suitable crop for cultivation. This platform facilitates the application of machine learning in real-world agricultural decision-making.

4.3.2 Technologies Used

The development of the website leveraged several technologies to ensure a seamless user experience:

- **Backend:** The backend was built using Flask, a lightweight and flexible Python web framework. Flask handles HTTP requests, processes input data, performs model inference, and returns prediction results.

- **Machine Learning Models:** The models were trained using various machine learning algorithms and were serialized using Python's pickle module for integration into the Flask application.
- **Data Scaling:** Two scalers, StandardScaler and MinMaxScaler, were employed to preprocess the input features, ensuring the data is appropriately scaled before being fed into the models.
- **Frontend:** The frontend was developed using HTML for structure, CSS for styling, and JavaScript for any necessary client-side interactions. The website's user interface allows for easy input of data and displays prediction results clearly.

4.3.3 Website Features

The website includes several key features:

Home Page: The "Home" page serves as the landing page of the website, providing an overview of the crop recommendation project and guiding users to the main functionalities.

Predict Page: The "Predict" page is the core of the website, where users can input relevant agricultural data. Users can also select the machine learning model they wish to use for the prediction. Upon submission, the website processes the input data and provides a recommendation on the best crop to cultivate based on the data.

Model Selection : The website allows users to choose from different machine learning models, enabling them to compare results or use a preferred model.

Prediction Display: The website displays the prediction result, indicating the best crop for cultivation given the provided environmental and soil data. If the data does not match any crop in the system's knowledge base, an appropriate message is displayed to inform the user.

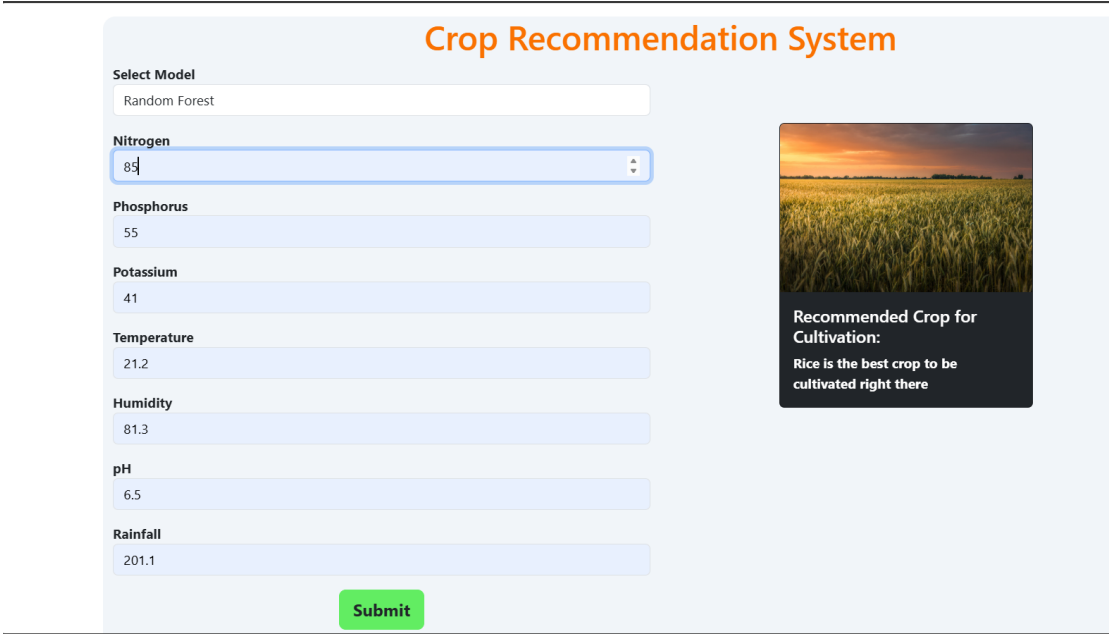


Figure 4.4: Crop recommendation using Certain values

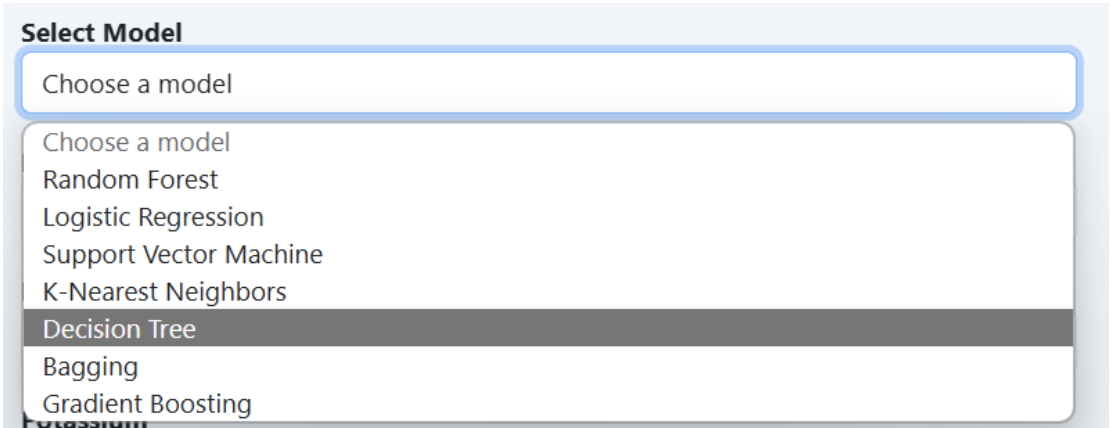


Figure 4.5: Selection of Model

CHAPTER 5

CONCLUSION

5.1 Summary

In conclusion, this project on Crop Recommendation using Machine Learning provides a valuable tool for optimizing agricultural practices by recommending the most suitable crops based on environmental factors. By analyzing key parameters such as soil pH, temperature, humidity, and nutrient levels, the system uses advanced machine learning models to deliver precise crop suggestions tailored to specific conditions.

The project incorporates various machine learning models, including ensemble techniques like Random Forest and Gradient Boosting, which are particularly effective at capturing the complex interactions between different environmental variables. The Random Forest model, in particular, demonstrated superior accuracy in predicting the best crops to cultivate.

Extensive data preprocessing, including feature scaling, encoding, and balancing, was carried out to prepare the dataset for model training, ensuring that the models performed effectively. Ultimately, this system represents a significant advancement in precision agriculture, aiding farmers in making informed decisions that can lead to improved crop yields and more sustainable farming practices.

5.2 Future Scope

- **Integration of Real-Time Data:** Incorporating real-time data from sensors and IoT devices can enhance the system's accuracy and responsiveness. By continuously updating environmental conditions, the system could provide more dynamic and timely crop recommendations, adapting to changing weather patterns and soil conditions.
- **Regional and Climate Adaptation:** Developing region-specific models by incorporating diverse datasets from various geographic locations can improve the

system's accuracy across different climates and soil types. Tailoring the recommendations to local conditions would enhance the system's applicability and usefulness in varied agricultural settings.

- **Advanced Machine Learning Techniques:** Exploring and integrating more advanced machine learning techniques, such as deep learning and reinforcement learning, could further refine the recommendations. These techniques could capture more complex patterns and interactions in the data, leading to more precise and effective crop predictions.

5.3 Limitations

- **Limited Dataset Size:** The dataset used in the project contains 2,200 instances, which, while comprehensive, may not capture all possible variations in environmental conditions across different regions. A larger and more diverse dataset could improve the model's generalizability.
- **Geographic and Environmental Generalization:** The model's accuracy may be limited when applied to regions with different climatic or soil conditions than those represented in the training data. Without region-specific data, the system might provide less accurate recommendations, highlighting the need for localization to different agricultural environments.
- **Dependence on Input Data Quality:** The accuracy of the crop recommendations heavily depends on the quality and accuracy of the input data, such as soil pH, temperature, and nutrient levels. Any errors or inconsistencies in this data can lead to incorrect predictions, highlighting the need for precise and reliable data collection methods.

5.4 Novelty

The novelty of this Crop Recommendation System lies in several key advancements that distinguish it from traditional approaches. First, the project significantly enhances the feature set by including crucial variables like nitrogen (N), phosphorus (P), and potassium (K) levels, along with other environmental factors. This comprehensive feature selection allows the model to make more accurate and context-sensitive predictions.

Moreover, the project incorporates robust data preprocessing techniques, such as Min-Max Scaling, standardization, and encoding, ensuring that the data is well-prepared and suitable for model training. Exploratory Data Analysis (EDA) was also employed to better understand the data, identify patterns, and guide feature engineering decisions.

The use of seven different machine learning models further adds to the novelty by providing a comparative analysis of their performance. This approach not only identifies the best-performing model but also offers users the flexibility to choose a model based on specific needs or preferences. Finally, the deployment of a user-friendly website where users can select a model and predict the optimal crop for their conditions adds a practical and accessible dimension to the project, making it a valuable tool for farmers and agricultural experts.

REFERENCES

- [1] Bondre, D. A. and Mahagaonkar, S.: 2022, Prediction of crop yield and fertilizer recommendation using machine learning algorithms., *Research Article* .
- [2] Garanayak, Mamata, G. S. S. M. and Jagadev., A. K.: 2021, Agricultural recommendation system for crops using different machine learning regression methods., *IGIGlobal* .
- [3] Reddy, D. Anantha, B. D. and Watekar, A.: 2019, Crop recommendation system to maximize crop yield in region using machine learning, *Research Article* .
- [4] Shilpa Mangesh Pande, D. P. K. R.: 2021, Crop recommender system using machine learning approach, *Proceedings of the Fifth International Conference on Computing Methodologies and Communication (ICCMC)*, IEEE, pp. 1–6.
- [5] Suresh, G., A. K.: 2021, Efficient crop yield recommendation system using machine learning for digital farming, *Research Article* .