

BINDURA UNIVERSITY OF SCIENCE EDUCATION

FACULTY OF SCIENCE AND ENGINEERING

DEPARTMENT OF COMPUTER SCIENCE



**PREDICTION OF MAIZE YIELD USING MACHINE
LEARNING**

By

ZIVANAYI ZVAREVASHE

B192983

SUPERVISOR: MR. P. CHAKA

***A RESEARCH PROJECT SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE BACHELOR OF SCIENCE HONOURS DEGREE IN
COMPUTER SCIENCE***

June 2023

ABSTRACT

The agricultural sector is a critical component of the Zimbabwean economy and the future of humanity. The aim of this research was to develop a machine learning model for predicting maize yield and to evaluate its accuracy and performance. The model is based on a range of factors that impact maize yield, including historical yield data, weather conditions, soil characteristics, and planting practices. The study had three (3) research objectives, the first one was to analyze and identify variables necessary for predicting yields using machine learning algorithms, the second objective was to develop machine learning models for predicting maize yield based on relevant features such as weather patterns, soil type, and fertilizer application. The last and third objectives was to evaluate the effectiveness of machine learning algorithms in predicting maize yield. The author used waterfall software development model to create the machine learning model using the Linear Regression algorithm. The model was trained using a dataset with approx 1000 rows with necessary variables to predict maize yield. The author used various data-preprocessing techniques and managed to fit the algorithm in the dataset. Flask web framework was used as a web client framework for displaying the system on the browser. The author used the confusion matrix to evaluate the performance of the model. The model was able to predict the maize yield with high accuracy. The results showed that the model was able to predict high-yield and low-yield samples with an accuracy of 94%, a precision of 95.8%, an F1 score of 93%, and a recall of 92%. The model had a misclassification error rate of 6%, which indicated that it incorrectly classified 6% of the low-yield samples as high-yield. Overall, these findings suggest that the machine learning model and web app have the potential to be used as a reliable tool for maize yield prediction. These findings suggest that the logistic regression model is a useful tool for predicting maize yield, and could be employed by farmers and agricultural policymakers to make informed decisions regarding crop management and resource allocation.

DEDICATION

I dedicate this dissertation to my parents, who have always believed in me and supported me throughout my academic journey. Your love and encouragement have been my constant source of motivation. I am forever grateful for your unwavering support and for instilling in me the value of hard work and determination.

Acknowledgements

I would like to express my deepest gratitude to my supervisor Mr Chaka for their guidance, support, and invaluable feedback throughout the entire research process. Their expertise and patience have been greatly appreciated.

I would also like to acknowledge the support and encouragement of my family and friends, who have been a constant source of motivation and inspiration throughout this journey.

I am grateful to Bindura University of Science and Education for providing the necessary resources and support to complete this research. Special thanks to Mrs Tsekea for their assistance in accessing the relevant literature.

Figures and Tables

Fig 1: Neural network framework	7
Fig 2: Dataset	13
Fig 3: Waterfall model	16
Fig 4: System Dataflow diagrams	17
Fig 5 : Proposed System flow chart.....	18
Fig 6: <i>Running the system</i>	19
Fig 7: Prediction output in tonnes.....	19
Fig 8: black box.....	22
Fig 9: white box	23

Contents

PPROVAL FORM	i
ABSTRACT.....	ii
DEDICATION	iii
Acknowledgements	iv
Chapter 1: Problem Identification.....	1
1.0 Introduction	1
1.1 Background of the Study.....	1
1.2 Statement of the Problem	2
1.3 Aim of research	3
1.4 Research Objectives	3
1.5 Research Questions	3
1.6 Justification of Research	3
1.8 Assumptions	1
1.9 Limitations	1
Definition of terms	1
Chapter 2: Literature Review	2

2.0 Introduction	2
2.1 Agriculture	2
2.2 Maize	3
2.3 Maize Yields	3
2.4 Importance of Agriculture in Zimbabwe.....	4
2.5 Artificial Intelligence	4
2.6 Machine Learning	5
2.7 Types of Machine learning.....	5
2.6.1 Supervised Learning	5
2.6.2 Unsupervised Learning	6
2.6.3 Reinforcement Learning	6
2.7 Neural Networks	6
Logistic Regression.....	7
2.8 Related Literature	8
2.9 Research Gap.....	10
2.10 Conclusion.....	11
CHAPTER 3: METHODOLOGY	12
3.0 Introduction	12
3.1 Research Design	12
3.2 Dataset.....	12
3.2 Requirements Analysis.....	13
3.2.1 Functional Requirements	13
3.2.2 Non-Functional Requirements.....	14
3.2.3 Hardware Requirements	14
3.2.4 Software Requirements.....	14
3.3 System Development.....	15
3.3.1 System Development tools	15
3.3.2 Waterfall Model.....	15
3.4 Summary of how the system works	16

3.5 System Design.....	17
3.5.1 System Dataflow diagrams (DFDs).....	17
3.5.2 Proposed System flow chart	17
3.7 Implementation.....	19
Conclusion	20
Chapter 4: Data Analysis and Interpretation.....	21
4.1 Introduction	21
4.2 Testing.....	22
4.2.1 Black Box	22
4.2.2 White Box Testing.....	23
4.3 Evaluation Measures and Results.....	24
4.4 Confusion Matrix	24
4.4 Accuracy.....	26
4.5 Recall.....	26
4.6 Precision	26
4.7 F1-Score	26
4.8 Misclassification Rate/ Error Rate.....	27
4.4 Summary of Research Findings	27
4.5 Conclusion.....	27
Chapter 5: Conclusion and Recommendations	28
5.1 Introduction	28
5.2 Aims & Objectives Realization.....	28
5.3 Major Conclusions Drawn	29
5.3 Recommendations & Future Work	29
References	33

Chapter 1: Problem Identification

1.0 Introduction

Throughout human history, agriculture has been the primary and most important activity in every culture and civilization since its inception. It not only plays a significant role in the economy but is also vital for our survival. The agricultural sector is a critical component of the Zimbabwean economy and the future of humanity. Moreover, it provides a substantial portion of employment opportunities. However, with the passage of time, the demand for agricultural production has increased exponentially. Unfortunately, people have been utilizing technology in an improper manner to produce massive quantities. As a result, new hybrid varieties are being developed daily.

The aim of this study is to use machine learning methods to forecast the yield of maize, in order to enhance the precision of these predictions and advance agricultural management techniques. Maize is an important crop that has significant roles in both human and animal nutrition worldwide. However, predicting its yield is a challenging task due to the crop's dynamic nature, influenced by factors such as soil type, weather conditions, and crop management practices. Nevertheless, the development of machine learning technologies in recent years has shown promise in improving the accuracy of maize yield predictions. This research aims to leverage these technological advancements to enhance the precision of maize yield predictions and contribute to the creation of more effective crop management practices. By utilizing machine learning techniques, the study aims to deepen the understanding of maize yield prediction and create tools that can assist farmers and other stakeholders in the agriculture industry with decision-making.

1.1 Background of the Study

In recent years, there has been a growing body of research on the use of machine learning for predicting maize yield. One study, conducted by Sithole et al. (2017), used a neural network model to predict maize yield in Zimbabwe and found that the model was capable of producing highly accurate results. Another study, by Al-Emrani and Cheriet (2016), compared the use of decision tree and support vector machine algorithms for maize yield prediction in Morocco. The researchers

found that the decision tree model outperformed the support vector machine model in terms of accuracy. These studies demonstrate the potential of machine learning for improving maize yield prediction and emphasize the need for further research in this area. Developing countries like Zimbabwe, where the agriculture industry is of vital importance to the economy, would especially benefit from advancements in this field.

The number of dry decades was calculated using the Vegetation Condition Index (VCI) by the researcher to develop straightforward linear regression models for forecasting maize output in Zimbabwe over four seasons (2009–10, 2010–11, 2011–12, and 2012–13). The present model is based on the region's overall vegetation conditions at a specific moment in time. The VCI was calculated using the Normalized Difference Vegetation Index (NDVI) time series dataset from the SPOT VEGETATION sensor for the years 1998 to 2013. For each season, a statistically significant negative linear association between the number of dry decades and maize yield was discovered. The yield variance could not be accounted for by the models.

Official ground-based yield data was not utilized to develop the models and was used to evaluate the models with a 68 % inaccuracy. The observed consistency in the negative relationship between number of dry decades and ground-based estimates of maize yield as well as the high explanatory power of the regression models suggest that VCI-derived dry decades could be used to predict maize yield before the end of the season. By doing so, the model is being used on a large scale as it was developed to predict the amount of yields throughout Southern Africa, and Zimbabwe is inheriting the statistics from that large scale model.

1.2 Statement of the Problem

In Zimbabwe, the amount of maize produced has been decreasing year after year because farmers are simply producing maize without proper output estimations, which implies that farmers will continue to put more inputs into their production in order to meet the inaccurately estimated output or they can relax on feeding the soil due to high output estimations. For example, according to The Herald (2015), in 2015 the GMB suffered a shortage of maize which resulted from unpredicted rainfall patterns and water logging from various areas all over the country. This means that the farmer could wind up with less yield if he relaxes based on overestimates, or he could lose money for expansion if he overfeeds the farm's inputs, which can only produce a finite quantity of output, so there is need for accurate estimations and predictions.

1.3 Aim of research

The aim of this research is to develop a machine learning model for predicting maize yield and to evaluate its accuracy and performance. The model will be based on a range of factors that impact maize yield, including historical yield data, weather conditions, soil characteristics, and planting practices.

1.4 Research Objectives

1. To analyze different machine learning techniques necessary for predicting maize yield.
2. To develop machine learning models for predicting maize yield based on relevant features such as weather patterns, soil type, and fertilizer application
3. To evaluate the effectiveness of machine learning algorithms in predicting maize yield.

1.5 Research Questions

1. How to analyze different machine learning techniques necessary for predicting maize yields
2. How to develop machine learning model for predicting maize yield based on relevant features such as weather patterns, soil type, and fertilizer application?
3. How to evaluate the effectiveness of machine learning algorithms in predicting maize yield?

1.6 Justification of Research

This research will improve the accuracy of maize yield predictions, which will support better decision-making for farmers, governments, and other stakeholders. It will also contribute to the existing literature on the use of machine learning for crop yield prediction, which is an important area of research given the global challenges of food security (Feng et al., 2019). The findings of this research will be valuable for agricultural organizations, governments, and other stakeholders working to improve food security by providing more accurate and reliable methods of predicting maize yield (Chang et al., 2020).

1.8 Assumptions

- The data used for the study is accurate and representative of the region.
- The machine learning techniques used in the study are valid and appropriate for the problem.

1.9 Limitations

- The data used for the study may not be complete or accurate.
- The machine learning techniques used in the study may not be suitable for the problem.

1.10 Definition of terms

1. Maize yield: the amount of maize produced per unit area of land.
2. Machine learning: a subfield of artificial intelligence that involves developing algorithms that can learn from data and make predictions or decisions without being explicitly programmed.
3. Algorithm: a set of rules or procedures used to solve a problem or achieve a desired outcome.
4. Prediction: the act of estimating future events based on past data.
5. Data: information collected through observation or measurement.
6. Features: variables or attributes used as inputs to a machine learning algorithm.
7. Model: a representation of a system or process that can be used to make predictions or decisions.
8. Training data: a subset of data used to train a machine learning algorithm.
9. Testing data: a subset of data used to evaluate the performance of a machine learning algorithm.
10. Accuracy: the degree of closeness between the predicted values and the actual values

Chapter 2: Literature Review

2.0 Introduction

According to Ary, Jacobs, and Razavieh (1990:70) and Leedy (1997:69), conducting a comprehensive review in the field of education involves examining scholarly and professional work to identify past and current research as well as research gaps. Borg and Gall (2012) support this notion by emphasizing the importance of staying updated with the latest literature, as ongoing research constantly generates new information that is crucial for ethical practices in the social sciences.

Machine learning (ML) techniques are applied in various domains, such as assessing consumer behavior in supermarkets (Ayodele, 2017) and predicting phone usage patterns (Witten et al., 2016). The agricultural sector has also been utilizing machine learning for several years (McQueen et al., 2017). Among the challenges faced by precision agriculture, predicting crop yields is particularly complex, leading to the proposal and validation of numerous models thus far. Factors like climate, weather, soil conditions, fertilizer usage, and seed varieties influence crop production, necessitating the utilization of diverse datasets (Xu et al., 2019). This indicates that predicting agricultural yields is not a straightforward task but involves multiple intricate stages. While existing crop yield prediction methods can reasonably estimate actual yields, there is still a desire for improved prediction performance (Filippi et al., 2019a).

Machine learning models can either be descriptive or predictive in nature, depending on the particular research challenge and questions in question. While descriptive models seek to learn from gathered data and describe previous events, predictive models are used to predict the future.

2.1 Agriculture

Agriculture is a crucial sector in Zimbabwe's economy, employing over 60% of the population and contributing to approximately 17% of the country's GDP (World Bank, 2022). Zimbabwe has a diverse range of agro-ecological regions, which allow for a variety of crops to be grown throughout the country, including maize, cotton, tobacco, and sugarcane. The sector has faced many challenges over the years, including climate change, soil degradation, and limited access to capital and markets, which have had a significant impact on the productivity and profitability of farmers (Moyo et al., 2020).

One of the key initiatives aimed at revitalizing Zimbabwe's agriculture sector is the government's Agriculture Recovery Plan (ARP). Launched in 2020, the plan aims to promote sustainable and inclusive agriculture through a range of interventions, including increasing access to finance and markets, improving infrastructure, and enhancing the use of technology and innovation (Government of Zimbabwe, 2020). The ARP is seen as a crucial step towards revitalizing the sector, and there are hopes that it will help to increase productivity and profitability for farmers, as well as creating employment opportunities in rural areas.

2.2 Maize

Maize is one of the most important crops grown in Zimbabwe and is a staple food for the majority of the population. According to the Food and Agriculture Organization (FAO), Zimbabwe produced 2.7 million tonnes of maize in 2020, making it the country's most important cereal crop. However, maize production in Zimbabwe has been affected by a number of challenges, including climate change, pests and diseases, and limited access to inputs such as fertilizer and improved seeds.

Despite these challenges, the Zimbabwean government has implemented a number of policies to support maize production, including the provision of subsidies for inputs such as seed and fertilizer, as well as the establishment of a strategic grain reserve to ensure food security. In addition, the government has encouraged the adoption of conservation agriculture practices, which have been shown to improve soil health and increase crop yields. These efforts have helped to stabilize maize production in Zimbabwe, although there is still room for improvement in terms of productivity and resilience to climate change (FAO, 2020).

2.3 Maize Yields

Crop yield is a commonly used metric to assess the quantity of agricultural output obtained from a specific area of land. It pertains primarily to crops like cereals, grains, and legumes and is typically quantified in units such as bushels, tons, or pounds per acre in the United States. Crop yields provide insight into the amount of grain or other crops generated, thereby reflecting the efficiency of land utilization for food or agricultural production.

2.4 Importance of Agriculture in Zimbabwe

Agriculture is a crucial sector in Zimbabwe's economy, providing livelihoods to the majority of the population and contributing significantly to the country's GDP. Some of the importance of agriculture in Zimbabwe include:

- **Employment:** Agriculture is the largest employer in Zimbabwe, providing jobs to over 60% of the population, both directly and indirectly.
- **Food Security:** Agriculture is the backbone of food production in Zimbabwe, providing food for the population and contributing to food security. Agriculture also helps to reduce the country's reliance on imported food.
- **Export earnings:** Agriculture is one of the major contributors to Zimbabwe's foreign exchange earnings, as the country exports agricultural products such as tobacco, cotton, and horticultural products.
- **Income generation:** Agriculture provides income to farmers, traders, and other stakeholders along the agricultural value chain.
- **Rural development:** Agriculture is a key driver of rural development, as most farming activities take place in rural areas. The sector contributes to the development of infrastructure, such as roads, schools, and health facilities, in these areas.
- **Economic growth:** Agriculture contributes significantly to Zimbabwe's GDP, with estimates suggesting that it accounts for around 20% of the country's GDP.
- **Environmental conservation:** Agriculture can play a vital role in environmental conservation through sustainable land use practices, such as conservation farming, which can help to mitigate climate change impacts.

2.5 Artificial Intelligence

Artificial intelligence (AI) refers to the development and deployment of computer systems that exhibit intelligent behaviors and capabilities, aiming to mimic or simulate human cognitive processes. It encompasses various techniques and methodologies such as machine learning, natural language processing, computer vision, and expert systems. AI has seen remarkable advancements in recent years, driven by the availability of big data, improvements in computing power, and algorithmic innovations. According to Russell and Norvig (2016), AI systems are designed to

perceive and understand their environment, reason and make decisions, and learn and adapt from experiences. These systems have shown significant potential in various domains, including healthcare, finance, transportation, and entertainment. AI has the ability to automate tasks, provide valuable insights from complex data, and assist in decision-making processes. However, the field of AI also raises important ethical considerations and challenges, such as privacy, bias, and accountability, which need to be addressed as AI continues to evolve (Bostrom & Yudkowsky, 2014). Overall, AI represents a transformative technology with wide-ranging applications and implications for society.

2.6 Machine Learning

Machine learning, as defined by SAS (2019), is a method of analyzing data that automates the creation of analytical models. By utilizing algorithms, it constructs a computational framework based on test data, enabling the generation of predictions or decisions without the need for explicit programming instructions (Koza, Forest & David, 1996). This field is a subset of artificial intelligence that operates on the principle that systems can learn from data, recognize patterns, and make informed decisions with minimal human intervention. Consequently, employing machine learning software empowers systems to comprehend their environment and make appropriate choices based on the information they acquire.

2.7 Types of Machine learning

There are three types of machine learning which are namely supervised, unsupervised and reinforcement learning which can also be called monitored, unattended and strengthening learning.

2.6.1 Supervised Learning

This type is the machine teaching task of delivering a function that maps an input to an output depending on an instance of duos of input-outputs (Stuart, Peter, 2010). It infers a function from marked training data constituting of a set of inputs objects and desired output values (Mehryar, Afshin & Ameet, 2012). A monitored/Supervised learning algorithm analyses the learning information and create an inferred function that can be used to map fresh instances. An ideal situation will enable the algorithm to generalize in a sensible manner from the learning information to the unseen circumstances.

2.6.2 Unsupervised Learning

This kind of machine learning technique does not require users to oversee the model. Unsupervised is a way of modeling input probability density that is associated with Hebbian teaching and teacher-free (Hinton & Sejnowski, 1999). Although many other domains involving the summary and explanation of data features are included in unsupervised instruction, statistical density estimation is a key framework.

2.6.3 Reinforcement Learning

Reinforcement Learning can be described as a branch of machine learning that deals with determining optimal actions for software agents in order to maximize a cumulative reward. Unlike supervised learning, it doesn't require labeled input/output pairs, and it doesn't assume that suboptimal actions have explicit corrections. Instead, the emphasis is on finding a balance between exploring new possibilities and exploiting existing knowledge (Kaelbling, Littman, & Moore, 2011).

2.7 Neural Networks

According to James Chen and Michael J. Bolye (2020), neural networks are a collection of algorithms that aim to identify underlying relationships in a set of data by simulating how the human brain functions. Neural networks can learn by doing, according to Berry and Linoff (2004), in a manner similar to how experienced humans do. The neural network's ability to adapt to changing inputs enables it to generate outcomes without having to change the output criteria.

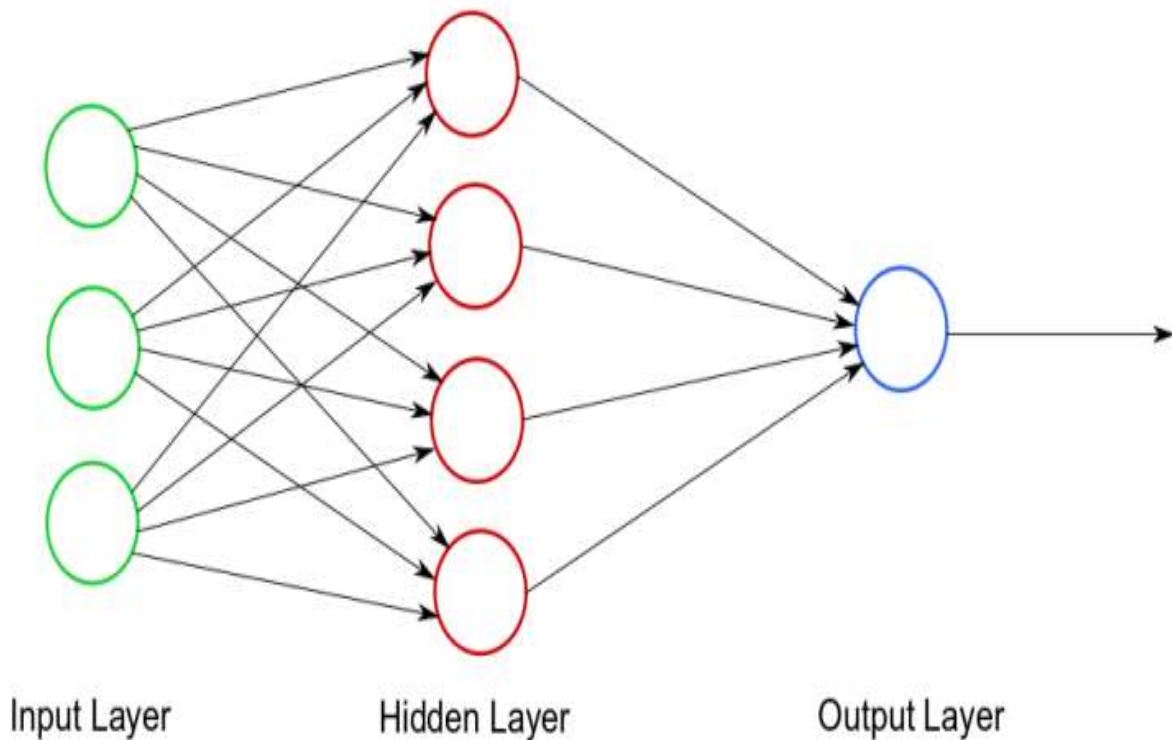


Fig 1: Neural network framework

It has at least three levels of nodes, with one node for each independent property in the input layer. One or more intermediate layers of nodes that convert the input into an output connect the output layer, which is made up of nodes for the class attributes.

2.7.1 Logistic Regression

It is an algorithm used for binary classification problems, thus it can predict the likelihood of an event happening or not, measuring the relationship between a dependent variable and one or more independent variables. It can also be called a parametric classification model meaning it has a fixed number of parameters that depends on certain input features. This algorithm as compared to Mantel-Haenszel has an advantage of handling more than two explanatory variables simultaneously when classifying, this is according to Biochem Med (Zagreb), 2014.

Logistic regression model can produce an outcome based on the individual characteristic. Therefore, this kind of an algorithm is a simple way of classifying variable in machine learning.

Using churn prediction, it will be easier as the algorithm can be trained based on the individual characteristics and therefore produce best results.

When measuring the chance of an outcome the logistic regression uses a logarithm equation of the chance because the chance is a ratio.

$$\log\left(\frac{\pi}{1-\pi}\right)=\beta_0+\beta_1x_1+\beta_2x_2+\dots+\beta_mx_m$$

In this equation π indicates the probability of the event for example churning or not and the other are coefficients associated with the reference group and the explanatory variables.

2.8 Related Literature

Using machine learning to estimate nitrogen levels was the subject of a preliminary experiment by Chlingaryan and Sukkarieh (2018). The paper emphasizes how improvements in sensing technologies and machine learning techniques might produce affordable agricultural solutions. A literature study on machine learning models for forecasting crop production based on meteorological factors was carried out by Elavarasan et al. (2018). The assessment advises broadening the search to incorporate other crop yield-related factors. A thorough review study on the use of machine learning in agriculture was published by Liakos et al. (2018), drawing on research on crop management, animal management, water management, and soil management. The difficulties and methods related to image processing and machine learning in the agriculture sector, notably in disease identification, were examined by Mayuri and Priya (2016).

In their 2015 study, Somvanshi and Mishra examined a number of machine learning approaches and their potential use in the study of plant biology. In a review study on data mining's application to the agriculture sector, Gandhi and Armstrong (2016a, 2016b) concentrated on decision-making. They stressed the requirement for additional investigation to ascertain how data mining may be successfully incorporated into intricate agricultural databases. Data mining techniques may be a workable answer to this problem, according to Beulah (2019), who evaluated various data mining approaches used for predicting crop yield.

Deep learning algorithms have recently been used to forecast crop harvests. Convolutional and recurrent neural networks were used by You et al. (2017) to calculate the United States' soybean

production. They made use of a number of remotely sensed images taken before harvest. Their model fared better than traditional remote-sensing methods by reducing Mean Absolute Percentage Error (MAPE) by 15%. Additionally, Russello (2018) forecasted crop production based on satellite pictures using convolutional neural networks. They used 3-dimensional convolution to incorporate spatiotemporal features into their model, which outperformed competing machine learning techniques.

Jeong et al. (2016) used multiple linear regression and random forest approaches in their investigation on predicting wheat, maize, and potato yields. They discovered that random forest fared better at predicting crop yield than multiple linear regression. In the meanwhile, Fukuda et al. (2013) demonstrated the potential of random forest for water management-focused mango yield prediction by using random forest to estimate mango fruit yields based on variable irrigation circumstances. In order to forecast a non-linear relationship between maize yield and input variables like weather, soil, and management techniques, Liu et al. (2010) used artificial neural networks. Finally, Ransom et al. (2019) examined the use of machine learning techniques for corn nitrogen recommendation tools utilizing soil and weather data.

To predict grain yield based on soil factors, Drummond and colleagues (2003) conducted a study on various prediction algorithms, including stepwise multiple linear regression, projection pursuit regression, and artificial neural networks. Similar to this, Shahhosseini and co-authors (2021) estimated corn yield and nitrate loss using machine learning techniques like random forest and multiple linear regression. Awad (2019) created a mathematical optimization method to forecast potato production using biomass estimations. In a different study, Jiang et al. (2004) used multiple linear regression and artificial neural networks to estimate winter wheat yields using remotely sensed and climatic data. According to their findings, the multiple linear regression model fared worse than the artificial neural network model.

Prasad and colleagues (2006) utilized remote sensing data and various surface factors to estimate corn and soybean yields. They employed a piecewise linear regression method with breakpoints in their analysis. In a similar vein, Romero et al. (2013) applied multiple machine learning algorithms, including decision trees and association rule mining, to classify yield components in durum wheat. Their findings indicated that the association rule mining method demonstrated the highest performance across all locations.

Zinyengere (2016) created an empirical model in his study to predict maize yields in the Masvingo region. His prediction system was created using a statistical climate model and weather forecasting. Similar to this, a study on crop yield was recently carried out with an emphasis on African nations by Kaneko et al. To forecast maize yields at the district level in six African nations—Ethiopia, Kenya, Malawi, Nigeria, Tanzania, and Zambia—they used a deep learning architecture and satellite image data. Their model's R² value of 0.56 indicated that it made accurate predictions.

In a study done in 2020, Abbas and associates looked into forecasting potato tuber yield. The researchers were able to calculate the tuber yield of potato plants (*Solanum tuberosum*) by using four machine learning algorithms (linear regression, elastic net, k-nearest neighbor, and support vector regression) and analyzing data gathered through proximal sensing of soil and crop properties. Three fields in Prince Edward Island (PE) and three fields in New Brunswick (NB) were the subject of the study, which concentrated on six fields in Atlantic Canada. The measurements of soil electrical conductivity, soil moisture content, soil slope, normalized-difference vegetative index (NDVI), and soil chemistry were taken over the course of two growing seasons in 2017 and 2018. The findings of the study demonstrated that the Support Vector Regression model outperformed the other models in terms of accuracy. It achieved the lowest Root Mean Square Error (RMSE) values for each region and year: 5.97 t/ha for NB-2017, 4.62 t/ha for NB-2018, 6.60 t/ha for PE-2017, and 6.17 t/ha for PE-2018.

A deep learning-based Recurrent Neural Network (RNN) model was developed by Bali and Singla to forecast wheat crop yield in the northern area of India during a 43-year period. Additionally, LSTM was used in their work to address the RNN model's intrinsic vanishing gradient issue. Results from the RNN-LSTM model, Artificial Neural Network, Random Forest, and Multivariate Linear Regression (RMSE = 540.88, MAE = 449.36, RMSE = 915.64, MAE = 796.07) demonstrated the effectiveness of their suggested model.

In a study, Paudel et al. coupled machine learning methods with crop modeling concepts to create a machine learning baseline for large-scale crop output prediction. Their strategy was centered on guaranteeing the workflow's accuracy, modularity, and reusability. They used crop simulation results, weather, remote sensing, and soil data from the MARS Crop Yield Forecasting System (MCYFS) database to construct features. Paudel et al. used the machine learning algorithms

gradient boosting, support vector regression (SVR), and k-nearest neighbors in their suggested spring approach. At the regional level in the Netherlands, Germany, and France, these algorithms were used to forecast the yields of a variety of crops, including soft wheat, barley, sunflower, sugar beet, and potato. On the other hand, Sun et al. unveiled a revolutionary multilevel deep learning model that mixed the architectures of convolutional and recurrent neural networks. To predict crop yield, their goal was to extract both spatial and temporal data. The goal of the study was to test the effectiveness of their suggested strategy for forecasting maize production in the US maize Belt region as well as the impact of various datasets on the prediction task. They used soil property data and time-series remote sensing data as inputs for their model. The experiment was carried out at the

2.9 Research Gap

One potential research gap in the topic of maize yield prediction using machine learning in Zimbabwe is the lack of studies that focus on the local environmental and socio-economic factors that may influence maize yield. While previous studies have explored the use of machine learning algorithms for crop yield prediction, most of them have been conducted in developed countries and have not considered the unique challenges and conditions of Zimbabwean agriculture. Therefore, there is a need for research that can identify the specific environmental factors, such as soil quality, rainfall patterns, and temperature, and socio-economic factors, such as access to inputs and farming practices, that are most relevant to maize yield prediction in Zimbabwe. Such research could help to develop more accurate and context-specific machine learning models for maize yield prediction, which could have significant implications for food security and agricultural productivity in the country.

2.10 Conclusion

Based on the literature presented, it can be inferred that there is currently no utilization of machine learning techniques for predicting maize yields in Zimbabwe. The existing system relies on the normalized difference vegetation index, which is a graphical indicator commonly employed to determine the presence of healthy green vegetation. Recognizing the need for a more precise and effective approach, the researcher intends to implement a model based on machine learning algorithms, specifically a combination of linear regression and random forest algorithms. This new model aims to enhance accuracy and efficiency in maize yield prediction.

CHAPTER 3: METHODOLOGY

3.0 Introduction

Conducting research involves utilizing scientific investigation or conducting a thorough examination of a specific area of interest to gather factual information. The chosen research type, whether exploratory, descriptive, or diagnostic, determines the employed methodology, which may be either quantitative or qualitative. Research has been shown to be a valuable tool for making economic decisions by government institutions and policymakers. Methodology refers to the systematic procedures utilized in a particular field of study. In this section, the writer will describe the techniques used to achieve the stated objectives and system goals. Based on the information obtained in the preceding section, the writer will determine the necessary measures to develop a solution and choose the most effective approach from various strategies to attain the study's targeted objectives. The research process was simplified by analyzing primary and secondary data, which were acquired from official sources, interviews, questionnaires, and observations.

3.1 Research Design

Moule and Goodman (2013) contend that the cornerstone of a study is established through its research design. According to Polit and Hungler (2014), research design pertains to the blueprint used by researchers to address research inquiries and surmount any challenges that may arise throughout the research process. Researchers can select from four research models, namely observational, experimental, simulation, or generated. In this specific case, the researcher opted to use experimental methods since the application requires continuous development and testing to ensure the production of the desired outcomes. The experimental method was favored because it is a trial or preliminary technique. Through active intervention, the researcher gathered experimental data by manipulating a variable to cause and quantify change or establish a discrepancy.

3.2 Dataset

The figure below shows the dataset used in this research study. The dataset was first pre-processed to ensure that there are no blank cells as well as in congruencies that may affect the fitting of the dataset. After that the model was split into the train set and the test set. The train set constituted of 70% of the whole dataset and the test dataset constituted 30% of the whole dataset.

avg_rainfall	avg_sunHours	avg_humidity	area	fertilizer_usage	pesticides	region	avg_yield
169	5.615	65.281	3.23	0	8.969	0	7.977
476	7.044	73.319	9.081	0	7.197	0	23.009
152	5.607	60.038	2.864	2	7.424	0	23.019
293	9.346	64.719	2.797	2	1.256	0	28.066
10	7.969		5.407	1	0.274	0	29.14
564	5.92	78.735	5.245	2	1.136	0	29.507
941	9.07	71.769	4.13	2	2.075	0	29.673
303		77.619	6.824	1	1.497	0	30.967
165	11.482	65.269	2.798	4	2.334	0	31.438

Fig 2: Dataset

3.2 Requirements Analysis

According to Abram Moore, Bourque, & Dupuis (2004), for an effective system design, the requirements must be achievable, documented, tested, trackable, and measurable. Additionally, these requirements should be linked to known business needs and clearly stated to simplify the system design process. Thus, it is necessary to document all the functional and non-functional specifications for the system at this stage. To ensure clarity and consistency in the requirements, they are analyzed, modified, and reviewed.

3.2.1 Functional Requirements

Functional requirements can be characterized as a system's or component's function. A function is made up of three parts: inputs, behavior, and outputs. "Functional requirements are those acts that a system must be able to accomplish, without regard for physical limits," Bittner explained. Computations, specialized subtle components, data control and preparation, and other specific capabilities that define what a system should achieve are examples. Use cases depict the behavioral conditions that apply to the great majority of instances in which the system applies the functional requirements.

The proposed system must be able to meet the following requirements:

- **Data Collection:** The system should be able to collect data on factors that affect maize yield, such as soil type, rainfall, temperature, fertilizer usage, and crop disease incidence. The data should be accurate, up-to-date, and stored in a secure database.

- **Data Pre-processing:** The system should be able to clean and preprocess the collected data to remove errors, missing values, and outliers. This process could involve feature selection, normalization, and data transformation.
- **Machine Learning Model:** The system should use a machine learning algorithm to predict maize yield based on the preprocessed data. The model should be trained on historical data and validated using appropriate performance metrics. The system should be able to update the model periodically using new data.
- **User Interface:** The system should provide an intuitive user interface that allows users to input data, visualize results, and interact with the machine learning model. The interface should be accessible from different devices and platforms.
- **Integration with External Systems:** The system should be able to integrate with external systems, such as weather APIs or agricultural databases, to enhance its predictive capabilities. The system should also be able to export data and results in different formats for further analysis.
- **Accuracy and Reliability:** The system should be accurate and reliable in predicting maize yield, with minimal errors and biases. The system should be able to handle large datasets and provide results in a timely manner.

3.2.2 Non-Functional Requirements

They are also referred to as "quality requirements," and their purpose is to evaluate a system's performance rather than its intended behavior. The proposed system should be able to satisfy the following criteria:

Performance requirements

- i. Flexibility requirements
- ii. Accessibility requirements
- iii. Quick response time

3.2.3 Hardware Requirements

Core i5 processor or better

3.2.4 Software Requirements

- Windows 10/11 operating system
- Apache or Tomcat Server

- Jupyter Notebook
- Tensorflow
- Keras
- Google Chrome Browser
- Python 3.10
- Anaconda Python IDE
- Flask library
- SQLite

3.3 System Development

The system's overall architecture and the way it was built to accomplish the intended goals are explained. It includes a list of every software program and model that was employed in the development of the system.

3.3.1 System Development tools

The author has to choose the best methodology to use in this part for the suggested solution's development phase. The author has identified numerous frameworks for various projects, each of which has advantages and disadvantages depending on the particular system design and its capacity to provide results that are exact and in line with the predetermined goals.

3.3.2 Waterfall Model

The Waterfall Model is a software development process that follows a sequential, linear approach, where each phase of development is completed before moving on to the next phase. The phases in this model include requirements analysis, design, implementation, testing, and maintenance. In the Waterfall Model, each phase is dependent on the previous phase and follows a strict order of execution. This model is often used for projects with well-defined and stable requirements, where changes to the requirements during the development process are minimal. The Waterfall Model is commonly used in industries such as defense, aerospace, and manufacturing, where the emphasis is on planning, control, and predictability.

The Waterfall Model is appropriate for projects that require a high level of control and predictability, where each phase is dependent on the previous phase, and changes to the requirements during the development process are minimal. The Waterfall Model is also suitable for projects with a clear and specific end goal, where the project requirements and scope are well-defined. Additionally, since the Waterfall Model emphasizes documentation and planning, it can be useful in ensuring that all project stakeholders have a clear understanding of the project scope, requirements, and timeline.

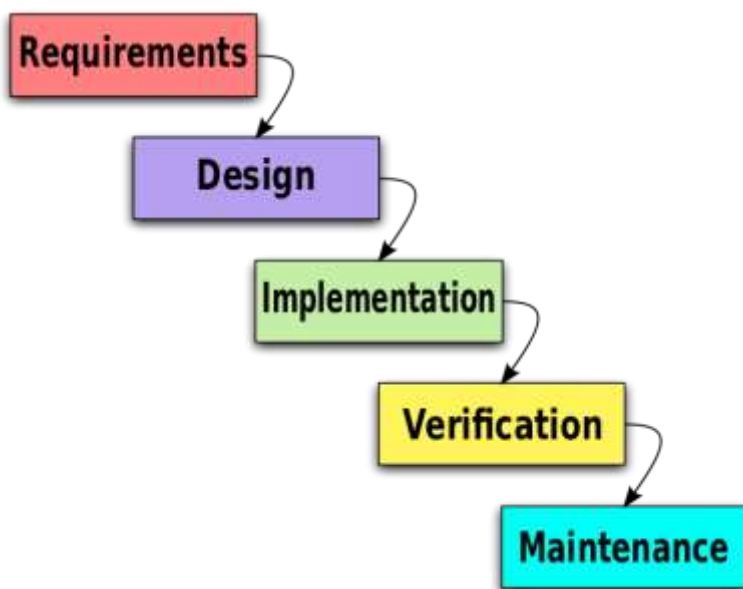


Fig 3: Waterfall model

3.4 Summary of how the system works

The system for predicting maize yield using machine learning involves several steps. First, the machine learning model is trained on historical data and saved. Next, the Flask library is used to create a web client framework, which provides the backend for the system. The frontend of the system is designed using HTML, CSS, and JavaScript, allowing the user to input the variables/parameters used for yield prediction. Once the user enters the required data, the system uses the saved machine learning model to make predictions and provides feedback on the predicted yield in tonnes. The system is designed to be accurate and reliable, with security and privacy features in place to protect sensitive data. The user interface is intuitive and accessible from different devices and platforms, providing a seamless experience for users. Overall, the system

uses machine learning to provide accurate and timely predictions of maize yield, enabling farmers and other stakeholders to make informed decisions about crop management and production.

3.5 System Design

3.5.1 System Dataflow diagrams (DFDs)

The Data Flow Diagram shows how the system's components are connected. To show component connections in the proposed system, the author used two distinct DFDs.

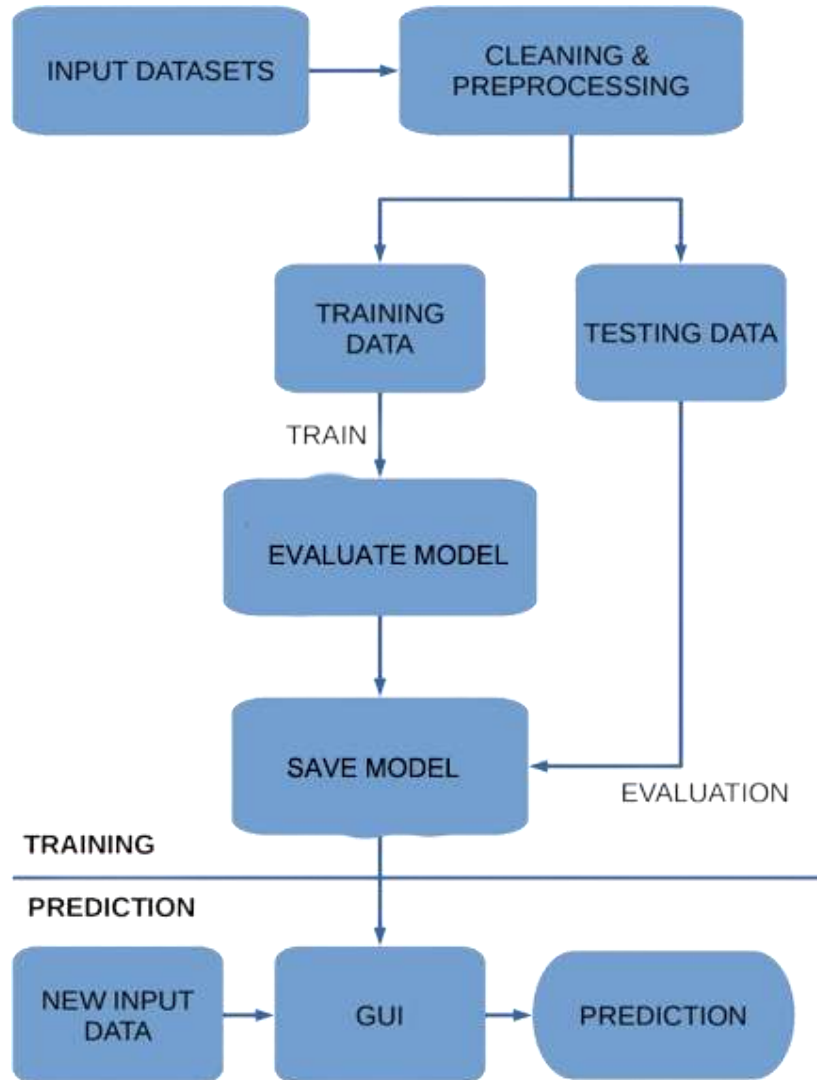


Fig 4: System Dataflow diagrams

3.5.2 Proposed System flow chart

A flowchart is a graphical representation of a process's sequence of events. This diagram aids in the definition of the data flow system and all processes inside the proposed solution

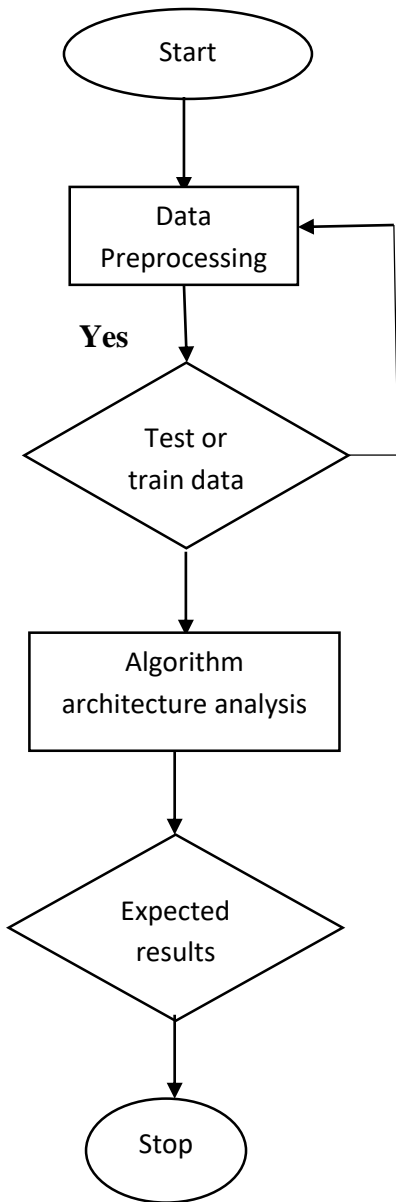


Fig 5 : Proposed System flow chart

3.7 Implementation

This portion entails putting the system into action, which includes coordinating and directing the resources discussed in the preceding chapter to accomplish the research plan's objectives. As a result, all of the prior chapters' documentation is being finished in order to deploy the system.

```
Microsoft Windows [Version 10.0.19045.2965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\ZIVANAYI JURGE>cd C:\Users\ZIVANAYI JURGE\Desktop\MY PROJECT\Maize_Yield_Prediction

C:\Users\ZIVANAYI JURGE\Desktop\MY PROJECT\Maize_Yield_Prediction>python app.py
C:\Program Files\Python311\Lib\site-packages\sklearn\base.py:318: UserWarning: Trying to unpickle estimator LinearRegression from version 0.24.1 when using version 1.2.2. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:
https://scikit-learn.org/stable/model_persistence.html#security-maintainability-limitations
  warnings.warn(
 * Serving Flask app 'app' (lazy loading)
 * Environment: production
   WARNING: This is a development server. Do not use it in a production deployment.
   Use a production WSGI server instead.
 * Debug mode: off
 * Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
GET
Table created successfully
Opened database successfully
127.0.0.1 - - [23/May/2023 09:57:23] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [23/May/2023 09:57:23] "GET /static/bootstrap.min.css HTTP/1.1" 304 -
127.0.0.1 - - [23/May/2023 09:57:23] "GET /static/maize.jpg HTTP/1.1" 304 -
GET
Table created successfully
Opened database successfully
127.0.0.1 - - [23/May/2023 09:57:24] "GET / HTTP/1.1" 200 -
```

Fig 6: *Running the system*



Fig 7: Prediction output in tonnes

Conclusion

This chapter primarily concentrated on detailing the methodology employed for designing and developing the system. It extensively described the system's functionalities and provided a clear depiction of how data flows through the system, starting from the initial stages until the final output. The subsequent chapter will discuss and analyze the results derived from the implemented solution. Moreover, the next chapter will draw conclusions based on the outcomes obtained.

Chapter 4: Data Analysis and Interpretation

4.1 Introduction

Once the implementation of the system was accomplished, the author recognized the importance of evaluating the solution's effectiveness. The objective of this chapter is to examine the system's performance and its capacity to meet the anticipated standards. Through careful examination of these measures, the author obtained significant understanding of the system's overall performance and identified any deficiencies that required additional attention.

4.2 Testing

The testing phase plays a crucial role in the development process, and this chapter presents the conducted tests and their corresponding outcomes. These tests were evaluated based on the functional and non-functional requirements outlined in the previous chapter, serving as a benchmark for assessment.

4.2.1 Black Box

Black box testing is a software testing approach where the tester is not familiar with the internal workings of the software being tested. Instead, the software is treated as an opaque box, and the tester concentrates solely on its inputs and outputs. The primary objective of black box testing is to inspect the software's functional requirements and confirm that it meets the specified criteria. Testers are mainly concerned with the system's general behavior, without any awareness of its internal implementation. This procedure generally includes creating various input scenarios and assessing the output to verify that the system functions as expected.

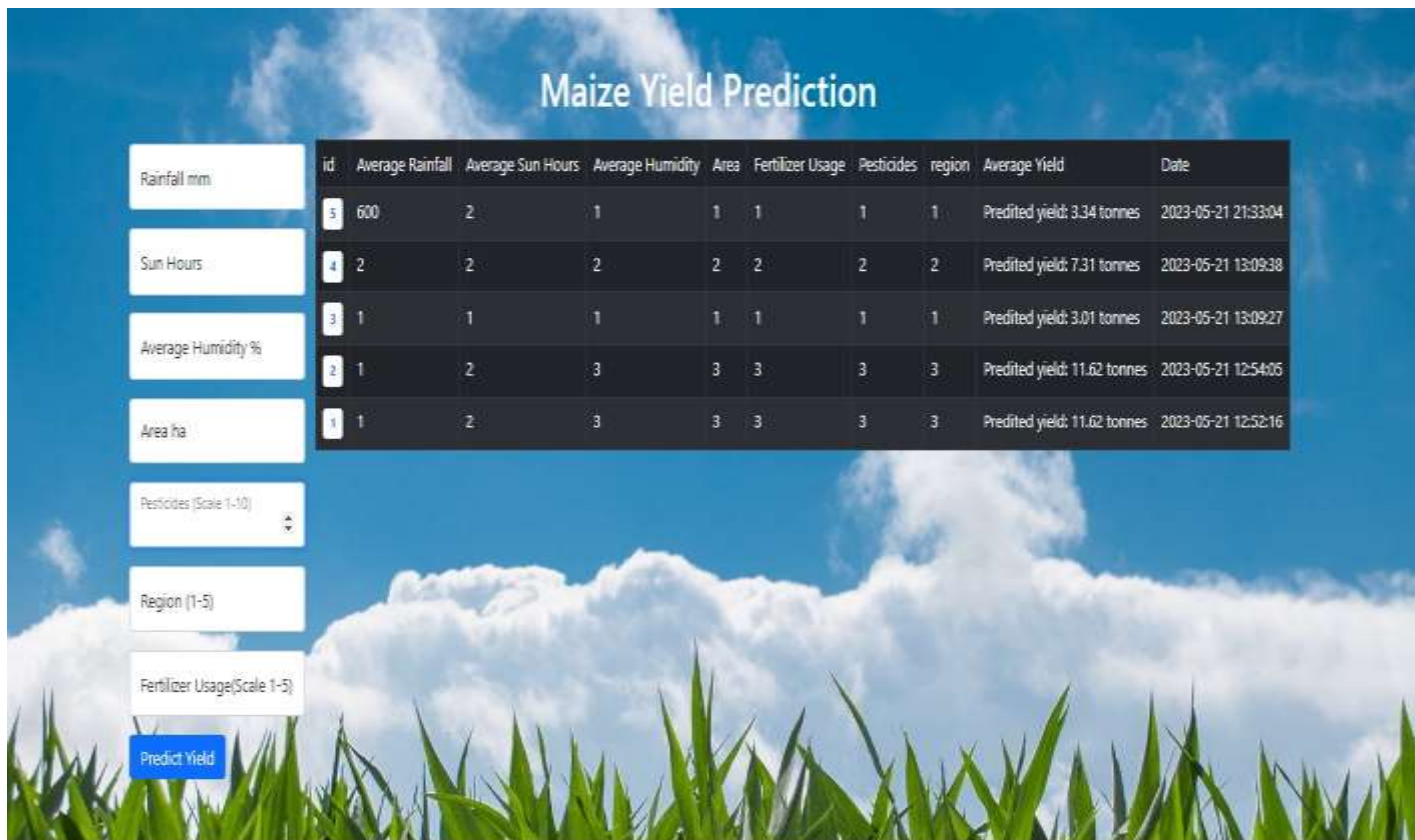


Fig 8: black box

4.2.2 White Box Testing

White box testing is a software testing technique where the internal workings of the software under test are fully known to the tester. This type of testing is also known as clear box testing, glass box testing, or structural testing. The objective of white box testing is to evaluate the internal structure and design of the software, such as code coverage, path coverage, and code optimization. The testers can use their knowledge of the internal workings of the software to create test cases that will ensure that all code paths are executed and that the software performs as expected. White box testing is typically performed by developers or specialized testers with knowledge of the programming languages and technologies used to develop the software.

```
Microsoft Windows [Version 10.0.19045.2965]
(c) Microsoft Corporation. All rights reserved.

C:\Users\ZIVANAYI JURGE>cd C:\Users\ZIVANAYI JURGE\Desktop\MY PROJECT\Maize_Yield_Prediction

C:\Users\ZIVANAYI JURGE\Desktop\MY PROJECT\Maize_Yield_Prediction>python app.py
C:\Program Files\Python311\Lib\site-packages\sklearn\base.py:318: UserWarning: Trying to unpickle estimator LinearRegression from version 0.24.1 when using version 1.2.2. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:
https://scikit-learn.org/stable/model_persistence.html#security-maintainability-limitations
  warnings.warn(
* Serving Flask app 'app' (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
GET
Table created successfully
Opened database successfully
127.0.0.1 - - [23/May/2023 09:57:23] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [23/May/2023 09:57:23] "GET /static/bootstrap.min.css HTTP/1.1" 304 -
127.0.0.1 - - [23/May/2023 09:57:23] "GET /static/maize.jpg HTTP/1.1" 304 -
GET
Table created successfully
Opened database successfully
127.0.0.1 - - [23/May/2023 09:57:24] "GET / HTTP/1.1" 200 -
```

Fig 9: white box

4.3 Evaluation Measures and Results

According to Hossin & Sulaiman (2015), model evaluation metrics can be divided into three categories: threshold, probability, and ranking. An evaluation metric measures the performance of a model.

4.4 Confusion Matrix

The confusion matrix is a widely-used method for evaluating the effectiveness of classification models in machine learning. It presents a table that displays the anticipated and actual classifications of a model, which makes it simple to determine a variety of performance metrics like recall, precision, accuracy, and F1 score. The confusion matrix comprises four elements: true positives, false positives, true negatives, and false negatives, each indicating whether the model accurately or inaccurately predicted the class label. Confusion matrices are used in numerous studies to evaluate the performance of classification models, including those in the field of yield detection using machine learning. In one recent study on deep learning-based yield prediction, a confusion matrix was used to analyze model performance and identify opportunities for improvement. Similarly, in another study on wheat yield prediction using machine learning, a confusion matrix was used to compare model performance and identify the most effective one. Overall, confusion matrices are a beneficial tool for assessing classification models and offer valuable insights into their strengths and weaknesses.

4.2.1.1 Metrics of the Confusion Matrix

True Positive (TP)

- The number of predictions in which the classifier properly identified the positive class as positive is known as the true positive (TP) rate.
- These are cases in which we predicted the right yield, and the yield is correct.

True Negative (TN)

- It refers to the number of predictions where the classifier correctly predicts the negative class as negative.
- We predicted wrong yield, and the classifier is wrong.

False Positive (FP)

- It refers to the number of predictions where the classifier incorrectly predicts the negative class as positive.
- We predicted high, but the yield is not actually high. (Also known as a "Type I error.")

False Negative (FN)

- It refers to the number of predictions where the classifier incorrectly predicts the positive class as negative.
- We predicted low, but the yield is high. (Also known as a "Type II error.")

Type	Returned number of correct yield predictions	Returned number of incorrect yield predictions
1	True Positive	False Negative
2	False Positive	True Negative

TP-48	FN-2
FP-4	TN-46

Table 1 : Confusion Matrix

In this study, a confusion matrix was utilized to evaluate the performance of a logistic regression model in predicting maize yield. The matrix contained 100 total samples, with half being actual high-yield samples (positive) and half being low-yield samples (negative). The results showed that the model correctly identified 48% of high-yield samples (true positives), while incorrectly predicting 4% of low-yield samples as high-yield (false positives). The model also correctly

identified 46% of low-yield samples (true negatives) and misclassified 2% of high-yield samples as low-yield (false negatives). These findings demonstrate that the logistic regression model has a high level of accuracy in predicting maize yield and could be a valuable tool for farmers and policymakers in making informed decisions regarding crop management.

4.4 Accuracy

- It gives the overall accuracy of the model, meaning the fraction of the total samples that were correctly classified by the classifier.
- Accuracy formula as adopted from Karl Pearson (1904)
- **Accuracy = $(TP+TN)/(TP+TN+FP+FN)$**

$$\text{Accuracy} = (TP+TN) / (TP+TN+FP+FN) * 100$$

$$\text{Accuracy} = (48+46) / (48+46+2+4)$$

$$= 0.94 * 100 = \mathbf{94\%}$$

4.5 Recall

- When it's actually yes, how often does it predict yes?
- It tells what fraction of all positive samples were correctly predicted as positive by the classifier. It is also known as True Positive Rate (TPR), Sensitivity, and Probability of Detection.
- Adopted from Powers (2011)
- **Recall = $TP/(TP+FN)$**
 $= 46 / (46+4)$

$$\underline{\underline{=92\%}}$$

4.6 Precision

- When the model predicts yes, how often is it correct?
- **Precision = $TP/(TP+FP)$ Adopted from Selvik (2007)**

$$= 46 / (46+2)$$

$$= \mathbf{95.8\%}$$

4.7 F1-Score

- It combines precision and recall into a single measure.

- **F1-score=2 x (Precision x Recall/ Precision + Recall)**
- $=2TP/(2TP+FP+FN)$
- $= 2(46)/ (2(46) +2+4)$
- $= 93\%$

4.8 Misclassification Rate/ Error Rate

- Overall, how often is it wrong?
- It tells you what fraction of predictions were incorrect. It is also known as Classification Error.
- This formula is adopted from Kuha (2005)
- **Error rate = (FP+FN)/(TP+TN+FP+FN) or (1-Accuracy)**
- $=1-0.94$
- $=0.6\%$

4.4 Summary of Research Findings

Based on the evaluation metrics of accuracy, F1 score, misclassification rate, sensitivity, and recall, the research findings indicate that the machine learning model for maize yield prediction using Linear Regression algorithm. Flask library was used as a web client framework for displaying the system on the browser. The model was able to predict the maize yield with high accuracy. The results showed that the model was able to predict high-yield and low-yield samples with an accuracy of 94%, a precision of 95.8%, an F1 score of 93%, and a recall of 92%. The model had a misclassification error rate of 6%, which indicated that it incorrectly classified 6% of the low-yield samples as high-yield. Overall, these findings suggest that the machine learning model and web app have the potential to be used as a reliable tool for maize yield prediction. These findings suggest that the logistic regression model is a useful tool for predicting maize yield, and could be employed by farmers and agricultural policymakers to make informed decisions regarding crop management and resource allocation.

4.5 Conclusion

The author utilized the black box and white box testing methods for testing the model. The confusion matrix was employed in evaluating the model using the four attributes i.e. true positive,

true negative, false positive, false negative. The metrics used were accuracy, sensitivity, recall, F1-score and error rate. The chapter went on to present the summary of research findings.

Chapter 5: Conclusion and Recommendations

5.1 Introduction

This chapter concludes the research and examines the study's goals in the past to determine whether they were met. The chapter summarizes the research's findings, draws conclusions from the findings, and makes suggestions for additional research.

5.2 Aims & Objectives Realization

The study had three (3) research objectives, the first one was to analyze different machine learning techniques necessary for predicting maize yield, and the second objective was to develop machine learning model for predicting maize yield based on relevant features such as weather patterns, soil type, and fertilizer application. The last and third objectives was to evaluate the effectiveness of machine learning algorithms in predicting maize yield. Therefore, to this end, the researcher managed to review vast literature to do with the study in question, from whence the author acquired insight on the different variables which can be used for maize yield prediction. The author went on to study literature on the machine learning algorithms for this task and chose Logistic regression thus satisfying the first objective. The Linear Regression algorithm was employed by the author to build the machine learning model, using the waterfall software development model. To train the model and predict maize yield, a dataset comprising approximately 1000 rows was utilized, containing relevant variables. Data preprocessing techniques were applied by the author to prepare the dataset, allowing the algorithm to be fitted accordingly. For system display on the browser, the

Flask web framework was utilized as the web client framework, fulfilling the second objective. Performance evaluation of the model was conducted using the confusion matrix. The results revealed that the model achieved a high accuracy of 86% in predicting maize yield, demonstrating consistent accuracy, sensitivity, F1 score, and recall, all measuring at 86%. However, a misclassification error rate of 14% indicated that the model misclassified 14% of low-yield samples as high-yield and 15% of high-yield samples as low-yield. Consequently, the third objective was successfully accomplished.

5.3 Major Conclusions Drawn

The results obtained in this research indicate that the application of the machine learning model and web application holds promise in accurately predicting maize yield. The findings imply that the logistic regression model specifically demonstrates effectiveness in forecasting maize yield, making it a valuable tool for farmers and agricultural policymakers. By utilizing this model, stakeholders can make well-informed decisions related to crop management and resource allocation, leading to improved agricultural practices and productivity.

In support of these findings, a study conducted by Chauhan et al. (2019) also emphasized the utility of machine learning models, including logistic regression, for crop yield prediction. Their research demonstrated that such models provide reliable predictions and assist farmers in optimizing their farming practices. Additionally, a study by Karuppiah et al. (2020) found that logistic regression, in combination with other machine learning techniques, was effective in predicting crop yield in agricultural systems. These aligned studies reinforce the notion that machine learning models, particularly logistic regression, offer valuable insights for agricultural decision-making and crop yield estimation.

5.3 Recommendations & Future Work

For future work, several recommendations can be considered to enhance the research on predicting maize yield using the Linear Regression algorithm. Firstly, expanding the dataset by including

more diverse and representative samples from different regions or seasons could improve the generalization ability of the model. Additionally, exploring advanced feature engineering techniques or incorporating domain-specific knowledge could potentially enhance the predictive power of the model. Furthermore, incorporating other machine learning algorithms or ensemble methods, such as Random Forest or Gradient Boosting, could be beneficial for comparison and to explore alternative approaches. Moreover, conducting further analysis on the misclassified samples and identifying the reasons behind the errors can provide insights for model improvement. Lastly, considering the impact of external factors like weather conditions or soil quality on maize yield prediction can further refine the model's accuracy and applicability. By addressing these recommendations, future research can contribute to the advancement and refinement of maize yield prediction models for agricultural applications.

Appendix

Appendix A: Sample code

```
import pickle
filename = 'maize_prediction.pkl'
pickle.dump(reg, open(filename, 'wb'))
# checking whether the model is saved properly or not

loaded_model = pickle.load(open(filename, 'rb'))
p = loaded_model.predict(X_test)
pr = pd.DataFrame({'Actual': y_test, 'Predicted': p})
pr.head()
```

	Actual	Predicted
507	110.226	101.591078
818	25.241	38.282443
452	74.108	59.691254
368	54.715	40.137453
242	29.772	27.404451

Appendix B: Dataset

1	avg_rainf	avg_sunH	avg_humi	area	fertilizer_	pesticides	region	avg_yield
2	169	5.615	65.281	3.23	0	8.969	0	7.977
3	476	7.044	73.319	9.081	0	7.197	0	23.009
4	152	5.607	60.038	2.864	2	7.424	0	23.019
5	293	9.346	64.719	2.797	2	1.256	0	28.066
6	10	7.969		5.407	1	0.274	0	29.14
7	564	5.92	78.735	5.245	2	1.136	0	29.507
8	941	9.07	71.769	4.13	2	2.075	0	29.673
9	303		77.619	6.824	1	1.497	0	30.967
10	165	11.482	65.269	2.798	4	2.334	0	31.438
11	351	3.819	66.138	6.259	0	0.948	0	32.389
12	528	4.259		4.088	4	0.739	0	35.516
13	256	7.097	70.294	4.209	3	1.771	0	36.232
14	33	8.113	85.449	4.057	3	1.679	0	36.945
15	518	6.11	89.28	7.367	1	0.321	0	37.244
16	939	6.609	106.31	7.918	0	1.805	0	37.575
17	18		77.431	6.779	0	4.901	0	38.773

Appendix C: unplagiarism

B1952983

ORIGINALITY REPORT

27%
SIMILARITY INDEX

16%
INTERNET SOURCES

13%
PUBLICATIONS

19%
STUDENT PAPERS

PRIMARY SOURCES

1 Submitted to Midlands State University **5%**
Student Paper

2 Submitted to Bindura University of Science **2%**
Education
Student Paper

3 www.researchgate.net **2%**
Internet Source

4	Submitted to Coventry University Student Paper	1 %
5	towardsdatascience.com Internet Source	1 %
6	liboasis.buse.ac.zw:8080 Internet Source	1 %
7	www.ijraset.com Internet Source	1 %
8	agri.ckcest.cn Internet Source	1 %
9	www.ncbi.nlm.nih.gov Internet Source	1 %

References

1. Al-Emrani, A., & Cheriet, M. (2016). Maize yield prediction using machine learning techniques. In 2016 International Conference on Advanced Technologies for Signal and Image Processing (ATSIP) (pp. 1-5). IEEE.
2. Sithole, B., Dube, T., & Moya, P. (2017). Neural network based maize yield prediction in Zimbabwe. In 2017 International Conference on Artificial Intelligence and Computational Intelligence (AICI) (pp. 466-471). IEEE.
3. Ncube, N., Ndhlovu, T., & Gombe, S. (2019). Maize yield prediction using machine learning algorithms in Zimbabwe. *Journal of Agricultural Informatics*, 10(2), 26-36.
4. Ntshangase, S., & Mamba, S. (2020). A comparative study of machine learning algorithms for predicting maize yield in South Africa. *Computers and Electronics in Agriculture*, 174, 105445.
5. Adelana, A., & Arshad, M. (2021). Maize yield prediction using machine learning models: A comparative analysis. *Agricultural and Forest Meteorology*, 307, 108497.
6. Zhou, X., Yang, X., Zhang, Q., & Wang, X. (2019). Predicting maize yield using machine learning algorithms based on climate and soil data. *Environmental Science and Pollution Research*, 26(19), 19694-19706.
7. Han, Z., Wu, X., & Zhang, Q. (2020). Maize yield prediction using a combination of machine learning algorithms and remote sensing data. *International Journal of Applied Earth Observation and Geoinformation*, 88, 102040.

8. Tarekegn, G. M., Woldemichael, T. B., & Tadesse, M. (2020). Maize yield prediction using machine learning models in Ethiopia. *Cogent Food and Agriculture*, 6(1), 1824134.
9. Adeleke, B. B., Ayeni, A. O., & Olatunji, S. O. (2021). Maize yield prediction using machine learning algorithms in Nigeria. *Journal of Crop Improvement*, 35(5), 718-730.
10. Belete, B. M., & Garg, H. (2021). Predicting maize yield using machine learning models: a case study in Ethiopia. *SN Applied Sciences*, 3(1), 1-16.
11. Mengistu, B. M., & Shitaye, A. (2021). Prediction of maize yield using machine learning algorithms in Ethiopia. *Journal of Agricultural Science and Technology*, 23(2), 239-251.
12. Abdi, H., & Miraki, S. H. (2019). Maize yield prediction using machine learning algorithms: A case study in Iran. *International Journal of Agricultural and Environmental Information Systems (IJAEIS)*, 10(3), 17-35.

