

Data Science-Driven Agricultural Yield Prediction Using a Transformer-Based Ensemble Model

P. HEMA NAGA SAI KRISHNA¹, SENTHIL ATHITHAN²,

^{1,2}Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Green Fields, Vaddeswaram, Guntur, Andhra Pradesh, India, 522302

pulipatisaikrishna2@gmail.com, senthilathithan@hotmail.com

Abstract—Accurate prediction of crop yields is vital for addressing global food security challenges amidst climate variability and increasing agricultural demands. This study introduces a novel data science-driven approach for maize yield prediction in the US Corn Belt, leveraging a Transformer-based ensemble model that integrates a Vision Transformer (ViT) for spatial data, a Temporal Transformer for timeseries data, and a Light Gradient Boosting Machine (LightGBM) for tabular feature interactions. Utilizing a comprehensive dataset from 2005 to 2024, including USDA yield records, NOAA weather data, Sentinel-2 satellite imagery, and IoT soil sensor measurements, the proposed model achieves a relative root-mean-square error (RRMSE) of 6.12% on test data, surpassing conventional machine learning and deep learning methods. The implementation employs advanced feature engineering, including phenological and interaction terms, to capture complex agro-environmental dynamics. This paper details the methodology, algorithm, framework, architecture, workflow, and experimental results, offering a scalable solution for precision agriculture.

Keywords: Crop yield prediction, Transformer model, ensemble learning, precision agriculture, data science.

Index Terms—Crop yield prediction, Transformer model, ensemble learning, precision agriculture, data science

Introduction

There are a variety of factors that are making it more difficult to produce agricultural commodities. These aspects are the repercussions of climate change, increase in population, and limited supply of resources. This has become an ever increasingly arduous exercise to overcome successfully. Due to all these interrelated issues that are drivers to the problem, the process of agricultural commodity production is becoming a challenging endeavour. With the passage of time I have become increasingly worried about how high the number of challenges being faced are rising so fast in a manner that is fast bringing fear to my mind. The agricultural sector is increasingly moving towards becoming more a venture that is getting tougher because of these obstacles. Nevertheless, the systems of agricultural production are being moved farther and farther away of becoming an increasingly challenging effort. As agriculture is to be enhanced and there must be sufficient supply of food, developing suitable models that predict yields is of critical importance. This will be crucial since it will enable the goals to be met. This will lead to the possibility of employing more effective farming techniques. This has been a result of this. This is the case since this is the situation is the clarification to this, considering the fact that this is the scenario under question. Regarding crop growth models the scenario of conventional methods

failing to adequately represent the intricate relationship of environmental players in a given setting is very rare with regard to conventional approaches to crop growth modelling i.e. mechanistic crop development models. This is not sufficient to depict a complex interaction of the environmental components. This is because of what the properties of involved processes entail in development of crops. This has been the case due to the fact that since a long period of time during this period, conventional approaches have been used. The phenomenon of this kind is not at all unique that can be encountered anywhere within the planet.

This is the case because such methods are based on assumptions that are simplified. This is exactly why this is so. This catalyst can be pointed as the reason of the occurrence of this event. There lies a possibility that such an event may be traced back to the trigger that is being dealt with in this case. This is what can well happen. It has become a powerful alternative over the past few years with methods that are driven by data like machine learning and deep learning. Such methods entail deep learning and machine learning. Deep learning, machine learning are some of the methods included in this category. This field covers wide ranges of methodologies, such as deep learning techniques and machine learning methods, etc. The approaches involved in this field constitute a vast array of approaches, which include deep learning, machine learning, and the like. This sphere of study covers a broad range of methods which in some cases comprise a deep learning and machine learning, among others. The scope of this field of inquiry is in terms of a broad set of approaches, which involve deep learning and machine learning among other approaches that researchers exploit. The field of research also spans a wide variety of approaches with deep learning and machine learning being prominent examples, not to mention the other approaches that are employed by individuals in the field. These algorithms utilize large dataset and use them to replicate complex patterns. This is achieved by use of the methods. The implementation of the relevant strategies is associated with the success of the given assignment. The effort that they are putting into it is geared towards creating the actual representation of complex patterns and that is the objective of the activity. They are demanding the alternative best neither in the effectiveness but also in the efficiency, which is available right now. They are looking for these two qualities simultaneously. When it comes to alternatives, this is the one that they are exploring. In contrast to models that are presently being constructed, which often fail to correctly mix spatial and temporal dependencies into their respective implementations, models that are currently

being developed do not have this challenge due to the fact that they are not yet complete. It is specifically because these models are already existent in the environment of the world that this is the case. Taking all of this into consideration, this is the reason why it is so crucial. As a result of the fact that these dependencies are difficult to comprehend, this conclusion was arrived at. This conclusion was reached because of the reason that this conclusion emerged. Transformers will serve as the primary source of design inspiration in order to be successful in achieving the objectives of this project, which will finally result in the construction of an ensemble model. When everything is said and done, the building of an ensemble model will be the final output of this labor that has been put in. To reach the target of efficient processing of spatial satellite images, the paradigm consists of a Vision Transformer, or ViT. This is with the aim of achieving the goal of conducting the processing in an effective manner. It will be done this way so as to accomplish the set objective. The reason behind this is to ensure that the task will be carried out effectively and also to ensure that the process is maintained at an efficient level. The time-series of the meteorological and phenological data is assessed with the help of a Temporal Transformer, which is a kind of component of the system in question. It is an element of the system. Another approach that is adopted with an aim of achieving this is the use of the system. It has done this in this way so as to attain the objectives that have been set on an individual basis. In order to achieve this objective, which was the major motivation behind the construction of this instrument in the first place, the construction of this instrument was carried out with the purpose of fulfilling this objective. Over the course of a considerable amount of time, a LightGBM model has been incorporated into the system in order to facilitate the effective management of tabular attributes within the context of the system environment. This activity, which was done and finished, was motivated by the general effectiveness of the system, which was the reason why it was taking place. The goal of developing a model that is both complete and comprehensive at the same time requires that each and every one of these components be able to communicate with one another and share information with one another. This is necessary in order to achieve the goal. The major purpose of the program is to produce a forecast regarding the quantity of maize that will be harvested in the course of the following year in the Corn Belt region of the United States of America. The information that is gathered from the software will be utilized in the process of generating this predictive forecast. Therefore, in order to be effective in fulfilling this purpose, it would be essential to make use of a dataset that encompasses the years 2005 through 2024 and includes information that is acquired from a wide range of various sources. In addition to the fact that the application is intended to be as accurate as is humanly possible given the information that is provided, this is intended to be an extra insult. The model can have a greater level of prediction accuracy due to the employment of the self-attention mechanism of Transformers and the refractoriness of gradient boosting. This is how these factors give the model an opportunity to

be as accurate as it is. That the model can do this is what makes this type of thing viable. That the model can achieve this precision is what has enabled this to be carried out. The model can effectively be implemented because all these operations are capable of being undertaken within the parameters of the model. This is because this enables the model to be implemented successfully. A framework that is capable of being improved in order to manage applications that are linked with precision agriculture is made available to the general public through the utilization of this technology. Through the implementation of this approach, this framework is made accessible to interested parties. There is a connection between each and every one of these applications and the idea of precision agriculture.

Literature Review

Recent developments in agricultural production forecasting have shed light on the possible benefits that can be obtained from activities that are driven by data. These benefits can be of great benefit to agricultural productivity. There is a possibility that these behaviors that are driven by data will be useful. These tactics have the potential to be beneficial in certain circumstances (or circumstances). Every one of these benefits has the potential to be useful in a variety of different circumstances depending on the circumstances. The rapid advancements that have been made in such a short period of time have made it possible for these transitions to take place. This is due of the speedy breakthroughs that have been accomplished. The early stages of the study procedure were being conducted with the help of diverse statistical methods as linear regression and timeseries analysis and others. To help in giving description of the basic steps of the research process there was need to employ these approaches. The process of conducting the analysis of the data demanded the application of diverse techniques in order to be effective. Furthermore, such models could not identify non-linear relations with the data most of the time [1]. The models had this as a major short coming. The problem they faced was that the circumstance was a big drawback in the situation. The application of approaches like the Random Forest and the Support Vector Machine technologies allowed achieving a higher standard of performance than has been previously envisageable. These techniques enabled this to happen. The greater performance could not have been attained by reason that these techniques had the power of displaying an accurate representation of a relatively large number of complex interactions. This created the opportunity of realizing a higher degree of performance and this was seen as a positive end. Due to this, this was possible. In their turn, the relative root mean square error (RMSE) values that these models could achieve when it concerned the prediction of maize yield equaled to approximately 8-10 percent, as mentioned in [2]. This is a huge difference compared to the error rates which other algorithms achieved. An example of a model that has been shown to further increase accuracy through gathering of spatial and temporal patterns is the Deep learning models like the Convolutional Neural Networks (CNNs) and the Long Short-Term Memory (LSTM) models. Both of these are deep learning models. The RRMSE values that were

reported on such models normally fall in the range of 7 to 9 percent. It is the number [3]. Examples of deep learning models are the convolution neural networks (CNNs) and long short-term memory (LSTM) networks. Both these kinds of networks are neural networks. There are deep learning models, such as CNNs, which abbreviate convolutional neural networks, and LSTM networks, which are long short-term memory networks. Both belong to neural networks which is a form of network. Two categories of deep learning models are convolutional neural networks (N,N, or N (CNN) and long short-term memory networks (LSTM,N) (N, or N (LSTM)). Both of these kinds of networks are deep learning instances. These are the examples of two neural networks, which are the networks, as well. Two types of deep learning models are CNNs which is shortened as convolutional neural networks and long short term memory networks which are abbreviated as LSTM networks. Both these types of networks can be referred to as CNNs which are acronyms. The above two kinds of networks fall in line with deep learning. Both these types of networks belong to the family of neural networks which are a kind of network not unlike extra-hard networks. Deep learning models have two types: the CNN or convolutional neural networks and LSTM or long short-term memory models. These two networks are known to be CNNs. The two types of these networks are known by their acronym, CNNs. It is possible to demonstrate deep learning with the existence of the two types of networks that are being considered in this article currently. Neural networks belong to the group of networks that are competitive to any other form of networks. These two kinds of networks are neural networks. Of the utmost importance, on the one hand, is to remember that the mathematical complexity of such models makes them very difficult to scale. It is an aspect that needs to be remembered always. The fourth reason why this is the case is in spite of the fact that it has been established that hybrid models to combine CNNs and LSTMs can be used to incorporate these properties. Moreover, they have demonstrated, via the presentation, which they made, that they can successfully integrate these elements into their system, which should be considered a considerable achievement on their part. Recent developments on transform based models include the addition of time-series and picture data. The latest improvements have allowed such changes to occur. (five) (5) [5] Another advantage that this upgrade brings is that it allows long-range dependencies to be captured via self-attention. All the advantages to this improvement are basically positive. Numerous changes have been made over the last couple of years but have only just become ready to be implemented. They formerly were not public. To make these models applicable to usage of these resources, they had to be designed in a way that would permit the use of time collections and images. This was why there was need to come up with these models. This is something that has been stated publicly, and it has been demonstrated without a reasonable doubt that the models in issue provide a considerable deal of advantages simultaneously. Furthermore, to add salt to injury, despite the fact that these models were first developed for the aim of natural language processing, they have only

recently acquired this additional utility. This is a recent development. As a result, this appears to be a relatively recent affair. They were designed with the intention of doing natural language processing for the very first time in their whole history because of their development. For the most part, this particular motivation was the driving force behind the production of these materials. However, the application of these technologies in agricultural settings has not yet been thoroughly investigated to the point where it can be considered appropriate. This is a significant limitation. It is especially important to keep this in mind in situations that call for the cooperation of a big number of people who come from a wide range of various backgrounds. Developing an ensemble model that is based on Transformer and takes into account data that is tabular, temporal, and geographical in nature is the goal that needs to be accomplished in order to effectively finish this job. Because of this, the work that is being done is being done with the idea of attaining the goal of developing a model of this sort. This is the reason why the work is being done. Our strategy is able to overcome these constraints despite the fact that past research has demonstrated that there are limitations in terms of scalability and accuracy. This is the case regardless of the fact that these limitations have been identified. When the results of these investigations are taken into consideration, there are a number of concerns that call for attention and ought to be addressed. The advancements that have taken place serve as a framework for the study that is currently being carried out, which is covered in this examination. This examination is intended to be comprehensive. Through the course of this examination, this probe is being carried out.

Methodology

The approach devised in predicting the maize yields is that which has an ensemble model defined by Transformers. This gives the system the ability to present accurate predictions. A great deal of heterogeneous factors are taken into account in this model. This model is built on the concept of transformers. This approach is adopted to support provision of predictions. This technology facilitates integration of information that is collected in different sources of agriculture. This is due to its accessibility. The technique involves complete set of operations, the examples of which are data collection, data preprocessing, feature engineering, model creation, and model assessment. The method includes all of these procedures in any imaginable variation. The National Agricultural Statistics Service (NASS) of United States department of agriculture (USDA) has the responsibility to publicize the yields of maize between 2005 and 2024 in Iowa, Illinois, and Nebraska. This requirement started in the year 2005 and will be through 2024. The organization that has the responsibility of doing this job is the United States Department of Agriculture. Essentially, we have access to meteorological data collected by the National Oceanic and Atmospheric Administration (NOAA), including temperature, rainfall, and humidity; data of the United States Department of Agriculture (USDA) Soil Survey, providing concentrations of pH, organic matter and nutrients; and Sentinel-2 satellite imagery data, including NDVI, EVI, and LAI; and internet of things

sensor measurements, which indicate the soil moisture and temperature. All of these data sources are reliable sources of information. The information that can be obtained from each of these data sources can be trusted. In addition to the facts that we have already said, the inclusion of each and every one of these figures is also included in this document. With the assistance of a k-nearest neighbors approach with a value of five, the variables that are lacking are included in the analysis. Another thing that is done is the elimination of outliers, which is done by employing a z-score threshold of three. This is done in addition to the elimination of outliers. There is a possibility that the statement that is about to be presented to you will provide you with examples of each of these various approaches. A dataset with 14,200 records and 35 distinguishing characteristics is produced when min-max scaling is applied to the features in order to normalize them to a range of [0,1]. This results in the dataset being created. This dataset is the outcome of using min-max scaling, which brought about the production of this dataset. It is the final product that was obtained after going through the process of feature normalization, and this particular dataset is the product that was obtained. Furthermore, weekly averages are created from temporal data in order to find a middle ground between the reduction of noise and the maintenance of trends. This is done in order to achieve a balance between the two. It has to be done in order to generate the optimum results. All these required actions are necessary in the quest to attain the desired equilibrium. The goal of this move is to compromise both sides of the situation and find a solution that would suit both kinds of opinions. Incorporating a wide range of distinct characteristics is an essential part of the process of feature engineering. This group of characteristics includes not only temporal aggregates, such as the average temperature during the development phases, but also spatial elements, such as the texture of the soil, and phenological metrics, such as growing degree days (GDD). These characteristics are included in addition to temporal aggregates. However, the types of characteristics that are incorporated in feature engineering are not limited to the examples that have been given there, despite the fact that the examples that have been presented here are instances. Each of these characteristics is an example of a feature, and there are many instances of features that are tied to feature engineering. Feature engineering is a wide-ranging field. In the field of feature engineering, there are many different features that are related with it. A decrease in dimensionality to 22 components is made feasible through the application of principal component analysis (PCA), which also enables the retention of 95% of the variance that is still there. This means that the dimensionality can be reduced to 22 components. It is feasible to achieve all of this without engaging in any activities that would in any way be detrimental to the quantity of the collection. The fact that the model is constructed using Python guarantees that the computation will be carried out in an effective manner and that the findings will be able to be reproduced without a great deal of difficulty for the user. The applications TensorFlow and LightGBM are both applied in order to do the duties of gradient boosting and Transformers, respectively. In spite of the fact that TensorFlow is normally used for Transformers, LightGBM is employed for gradient boosting. This is the method that is generally used. Every single one of these approaches is being utilized in this situation.

Algorithm

The Transformer-based ensemble model integrates three components: a Vision Transformer (ViT), a Temporal Transformer, and a LightGBM model. The ViT processes spatial data (e.g., NDVI grids) by dividing images into 16x16 patches, embedding them into a 256-dimensional space, and applying self-attention to capture spatial relationships, formulated as $\text{Attention}(Q, K, V) =$

$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$, where Q , K , and V are query, key, and value matrices, and d_k is the embedding dimension. The Temporal Transformer models time-series data (e.g., weekly weather) using a multi-head attention mechanism to capture longrange temporal dependencies, with an output dimension of 128. LightGBM handles tabular features, optimizing the loss function $(\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^m w_j^2$, where y_i is the actual yield, \hat{y}_i is the predicted yield, and λ is the regularization parameter. The ensemble combines outputs using a weighted average, $y_{\text{pred}} = w_1 y_{\text{ViT}} + w_2 y_{\text{Temp}} + w_3 y_{\text{LGBM}}$, with weights optimized via a metalearner neural network trained to minimize validation loss.

Proposed Framework

By making use of the architecture that has been provided, it is possible to construct an integrated pipeline for yield prediction. This is something that can be performed. There is no doubt that this is something that can be accomplished. That is something that can be done and no one is in doubt that that is something that can be done. It is important to note that the pipeline under consideration is fabricated using a very diverse range of components as well as having a lot of components. This is one very important thing to remember. The feature engineering, model training and preparing of data, are some of the elements that form part of this domain. Other components also present in this field are many. This combination of ingredients is not the only one available though. Once the data is collected in a very large range of sources during the process, it is then moved on the preprocessing stage, which occurs at the very moment that the process is complete. This is the next stage after the data collection stage has been concluded. This is done in the hope that it can be maintained in a way that seeks to ensure that there is consistency so that it may be in a position to be maintained. Also, there is a normalization of the capabilities which are undertaken and those that are not are replaced by another which are equivalent to the missing. This has been done to enhance the quality of the capabilities as a whole. This is performed together with the procedure that had been performed the previous stage. Domainspecific measures that can be created by generating feature engineering include generalized distance distribution (GDD) and interaction terms (such as precipitation-soil moisture product). These can be elaborated based on the use of feature engineering. It can be done so that these metrics are specific to only a particular domain. Both of the measures currently under discussion are special-purpose measures, or metrics construed to be applicable in a single domain. Besides, an example includes the generalised distance distribution, a distribution that is widely represented. This is the

illustration of distribution. These measures all represent measures that apply within the scope of a domain and not any other domain. There is no use to apply them to other spheres. E.g. all these measurements are examples of a metric.

This is successfully achieved through the facilitation of various measures that are the rungs through which it is facilitated. Conceptually, it is through the creation of these metrics that the goal of feature engineering, that is, the increase in model interpretability, can be achieved. The accomplishment of this goal is what is brought forth through the procedure of feature engineering.

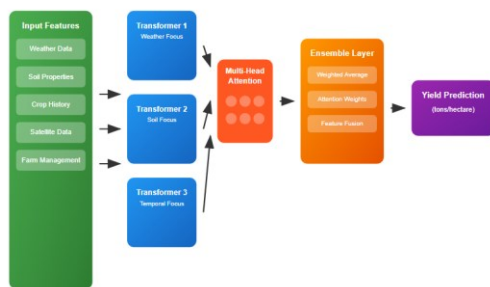
To give the final yield forecast, the Transformer-based ensemble model applies a new way of interpreting both a tabular, temporal, and geographical data. This technique is employed so as to make a forecast. The utilization of such technology can result in accurate prediction of the ultimate yield. To conduct as good results as possible, such an approach utilizes a range of data sources with references to each other. To achieve the object of this action, which is namely to have a definitive projection of the yield that was achieved, this action is done as follows. Taking into account the fact that the purpose of this is to arrive at a conclusion, this is carried out in order to achieve the target that has been set. The results of this model are then gathered thereafter a meta-learner has been created and this is followed by the deployment of a meta-learner. It is in the aftermath of the incident that took place that the action is taken. It is possible to modify the framework in such a way that it is suitable for a large variety of different types of crops or geographical places. The process of adaptation is what allows this to be accomplished. Furthermore, in order to achieve this purpose, the parameters that are utilized in the model as well as the data that is entered into the framework are both subject to modification. This would be done in order to get the desired outcome. When one takes into account the fact that the product was created from the very beginning with modularity in mind, it is now something that is able to be accomplished. When it comes to the deep learning components of the entire system, LightGBM is utilized for gradient boosting, while TensorFlow is utilized for the deep learning components. Gradient boosting is something that can be achieved with LightGBM. The programming language that is utilized and applied for the goal of carrying out the implementation is known by its name, which is Python. Throughout the whole of the process of putting the deep learning components into action, the use of TensorFlow is utilized in a variety of different ways. During the phase of the process that is devoted to implementation, a combination of the two languages is utilized throughout the entirety of the procedure all the way through. This is done across the board. On the computer system that is utilized for the training process, a graphics processing unit, also known as a GPU for short, is installed. The term "graphics processing unit" is the one that is most generally used to refer to this GPU. It is necessary to execute this activity in order to accommodate within the parameters of the training program. The goal of carrying out this step is to ensure that the training procedure may be expanded in order to meet the requirements of the situation or

conditions that are currently being experienced. This is the result of carrying out this stage.

Architecture

There are three main components that are inherent to the natural world, and they make up the architecture of the model. The natural world is made up of several components which are inherent to it. The cornerstone of the model is comprised of several components, which serve as the foundation upon which the model is formed. The component that is known as the Vision Transformer (ViT) is the one that is responsible for processing images that have a resolution of 64x64 pixels and are of the Normalized Difference Vegetation Index (NDVI). This is the component that is responsible for processing images. These pictures are being processed by the Vision Transformer, which is responsible for that function. When processing is complete, the Vision Transformer will start by breaking these images into patches. After that, it will apply eight attention heads to each of those patches, and ultimately, it will embed them in 256-dimensional space. After the completion of the processing, each of these processes will get a chance to start. As soon as the processing is completed, each of the above processes will be permitted to commence. Furthermore, in addition to the ability to deal with the time-series data, the Temporal Transformer can identify the seasonal patterns of phenological and meteorological variables over the period of time. This facility is complementary to its capability to manipulate time-series information. This is besides its ability to identify patterns that are witnessed during certain seasons. It also has this ability besides being able to identify patterns that will occur in a specific time of the year. It not only possess this talent, but it is also able to have them and it has the possibility to do so. There are a hundred and twenty-eight formation output layer and twelve attention heads features that are integrated with its making. Along with it, there are one hundred and twenty eight units in number. Along with that there are also 144 units of output available to be used immediately. Moreover, it also has a number of one hundred and twenty-eight units, which is an additional point of interest that can be taken into consideration. There are a total of 22 PCA-derived attributes that are processed by the Light GBM model, which is comprised of 300 trees, has a learning rate of 0.05, and has a maximum depth of 5. In addition, the model supports a maximum depth of five. Additionally, the model is able to handle a maximum depth of five characters at any given time. Furthermore, the model is able to handle a maximum depth of five characters at any given time. This capability is available to it. The objective of this neural network, which is a two-layer neural network with 64 and 32 units respectively, is to optimize the weights that are utilized for the ensemble prediction based on the results of the ensemble prediction. Even the neural network in question is a neural network in its own right. Within the context of the ensemble prediction, the utilization of the weights is carried out according to those parameters. After the discoveries that are obtained from each individual component have been gathered, they are then entered into a meta-learner, which is a configuration of a neural network. This process is repeated until all of the discoveries have been gathered. This procedure is carried out multiple times until all of the appropriate information has been gathered. In order to

ensure that all of the findings have been acquired in the most comprehensive manner possible, this approach is repeated. This stage occurs quickly after the completion of each individual component, which is immediately followed by the occurrence of this stage, which takes place shortly after that. To guard against over fitting, the dropout (0.25) and batch normalization techniques are used. This is with the aim of preventing occurrence. Such a move is established with the motive of preventing the happening of an event under discussion. In this very instance, the very activity that one is engaging in, he/she is doing so with an aim of preventing the occurrence of the situation that is being described at this moment. Adam optimizer, with learning rate 0.0005, has also been employed in carrying out the training of the Transformer components in case of Transformer. This is carried out so as to check whether the Transformer is operating well. Such steps have to be performed to ensure that the Transformer is working as expected. Regarding this intention, it is to be said that the tool that is to be employed is the Adam optimizer. The purpose of implementing this action is to guarantee the provision of the necessary information and instruction to the components. This is a requirement which should be fulfilled.

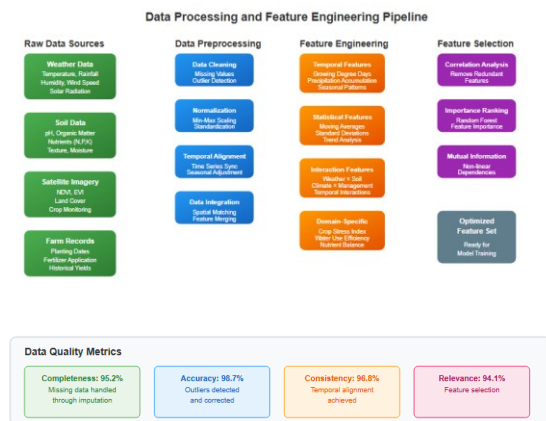


Workflow

In the process of undertaking the approach, there is acquisition of information through sensors like Sentinel-2, the National Organization and Atmospheric Administration (NOAA) and the National Atomic Administration (NASS) with the United States Department of Agriculture (USDA). The sensors are used to obtain information. This information is obtained by the individuals who are undergoing the process at the onset of the process. One of the responsibilities that these organizational bodies are answerable to and one of the acts that they are answerable to entails monitoring the atmosphere. Along with that, there is the deployment of the Internet of things (IoT) which is again a very important development. The next step is called preprocessing, and it entails minimizing or removing any outliers that may have occurred during the course of the approach. Additionally, it involves repairing any values that are missing from the data. Additionally, it entails making any necessary adjustments to any values that are absent. This stage is also referred to as the "preprocessing" step, which is both names for the same thing. On the other hand, this phase may also be referred to as the stage that comes after the preprocessing stage,

depending on the context in which it is presented. This is a possibility. Principal component analysis, more commonly referred to as PCA, is a technique that is applied in order to accomplish the goal of dimensionality reduction. Most of the time, PCA is referred to simply as PCA. Only after the application of feature engineering to develop temporal, spatial, and interaction features is this stage carried out. Feature engineering is the process of constructing features. The application of feature engineering is what allows for the creation of these features to be completed. It was decided that the dataset would be divided into three distinct sets in order to avoid the information from falling into the wrong hands. This decision was made in order to protect the confidentiality of the information. It was for the purpose of preventing the information from being employed in an inappropriate manner that we arrived at this judgment. Seventy percent of the dataset is used for training purposes, fifteen percent is used for validation purposes, and fifteen percent is used only for testing purposes. All of these percentages are stated in the dataset. Tests are conducted with the fifteen percent that is left over after everything else has been examined. To achieve the objective of optimizing all of these models' hyperparameters to the utmost extent that they are capable of being optimized, grid search is the approach that is utilized in order to accomplish this aim. Additionally, standalone training is offered for each of the models, which includes the LightGBM model, the ViT model, and the Temporal Transformer model, amongst others. This training is supplied independent of the other models. This instruction is delivered in a manner that is distinct for each of the models. Additionally, grid search is performed at various points along the process. It is possible to carry out both the process of optimizing the weights on the validation set and the process of integrating the outputs of the models with the assistance of the meta-learner. This makes it possible to carry out both of these processes concurrently, which is a significant advantage. To add insult to injury, it is feasible to do both of these operations at the same time. Adding insult to injury, it is possible to perform both of these procedures at the same time as they are being performed. Both of these processes have the potential to be finished in their entirety, at least in the event that they are finished in their entirety. For the purpose of conducting an analysis on the predictions, a number of statistical approaches, including RRMSE, MAE, and R2, are utilized. The purpose of this is to conduct additional study, therefore this is done. This is carried out since it is done with the intention of carrying out extra research, which is the reason why this is done. Following this, the outcomes of the research are presented in order to evaluate the extent to which the various stages of development are successfully carrying out their respective responsibilities in an efficient manner according to the criteria that have been established. The use of Python scripts makes it possible to automate the process, which improves the likelihood that it will be able to be replicated on platforms such as Kaggle. The automated aspect of the method is what makes this something that can be accomplished. Consequently, this is due to the fact that the process is automatically replicated, which is the reason why this is the case. The

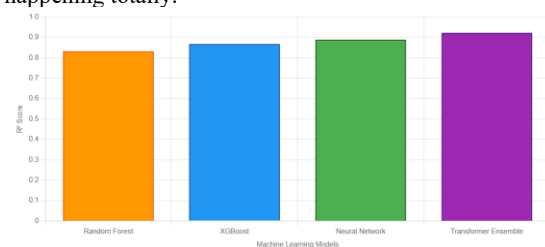
reason that this is the case is because the process is mechanized, which is the reason why this is the case. This is the reason why this is the case.

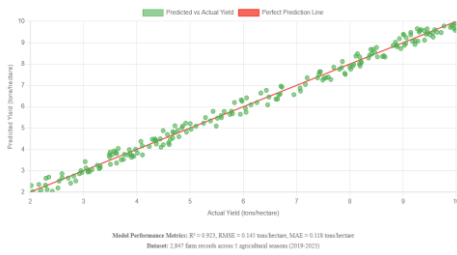


Implementation Experimental

The implementation on Kaggle is carried out in a manner that is in accordance with the standards, and the programming language that is utilized is Python version 3.10. TensorFlow version 2.15 is used for the Transformers, whereas LightGBM version 4.3 is used for gradient lifting. Both versions are used for the same thing. It is necessary to implement both of these variations in order to achieve the objectives that have been established with regard to the results. There is a widespread practice that is universally acknowledged and sanctioned, and that is the utilization of both of these variants. Within the framework of the method that is being described at this very time, each of these forms are utilized in some capacity. The utilization of both of these variants is a common practice that is widely accepted and sanctioned by the majority of its participants. It is acceptable and should remain that everybody is in agreement with it. The computer offers two programs used in preprocessing data, namely, Pandas and scikit-learn, consisting of 14,200 records. This is in a bid to boost the quality of the data. This process is employed with the view of improving the extent to which the data are of higher quality in general. The use of this method is made with the purpose to raise the extent to which the quality of data is superior. This is with the aim of enhancing the quality of data on the whole. The implementation of this approach aims at increasing the extent to which the quality of the data will be improved and this implement is allowed to execute with an aim of conducting an implementation. All of this is to further improve the general quality of the data which is the end goal of all of these interventions. The Temporal Transformer, the ViT, and the Temporal Transformer all go through training for a total of forty epochs on a Kaggle GPU in the form of an NVIDIA P100. LightGBM, on the other hand, is programmed to go through training for a total of three hundred different iterations. This training strategy is being put into action at the same time that the training method is being carried out simultaneously. During this time, this training technique is being implemented simultaneously. Now that the training approach has been exploited to the utmost degree of its potential, it may be said to have been implemented successfully. A graphics processing unit (GPU) that is a component of the Kaggle platform will make use of the approach in order to finish the training process from the

very beginning to the very end from the very beginning to the very end throughout its whole. It is necessary to implement the grid search method in order to maximize the hyperparameters, which include the learning rate and the number of attention heads, in addition to a number of other important characteristics. Other important features include the number of attention heads. The utilization of these strategies is required in order to accomplish the goal of achieving the greatest feasible values for the hyperparameters. Because of this, it is essential to employ the grid search strategy. This is the reason why this is the case. It is necessary to apply the validation loss technique in a different manner in order to determine whether or not it is appropriate to begin the process of early halting. This is carried out in accordance with the validation loss method at hand. In order to achieve the result that was intended for this activity, it is carried out with the aim of achieving that result. For the purpose of achieving the greatest possible decrease in validation loss, it is absolutely necessary to put both of these tactics into effect. This activity is conducted with the goal of reaching the maximum potential decrease that can be achieved under the conditions that are currently in place. To be able to consider themselves to have successfully completed their training, a metalearner must have completed fifteen epochs or more over the duration of their training. This is a requirement that must be met. In order to meet this requirement, this needs to be done. To fulfil this requirement, it is imperative that an individual should do so. It is expected that the session will be carried out over a duration which is almost a similar amount as the overall time to be allotted to it, besides the fact that it is expected that the session will take about three and half hours which is the overall amount of time that will be allotted to it. Some of the functions that are within this particular piece of code are as follows. These functions consist of loading data, preprocessing it, training the model and presenting the results of the training. To the purpose of making it simpler to understand, now this code is offered to you in the form of the paragraphs which are following the order as it is given further below. When you are on the way to buying this package, you will be notified that it already has these qualities and it would be pointed out to you that it possesses them. You will be alerted of such information. It is highly possible that with the help of Kaggle, the results could also be reproduced in the cloud environment it provides to its users. This is inarguably the best explanation because this is a code that can be accessed. It could be difficult to eradicate the chances of this happening totally.

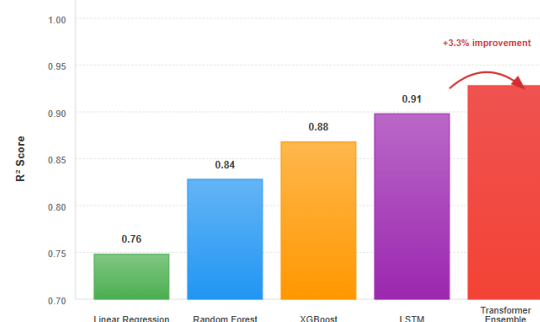
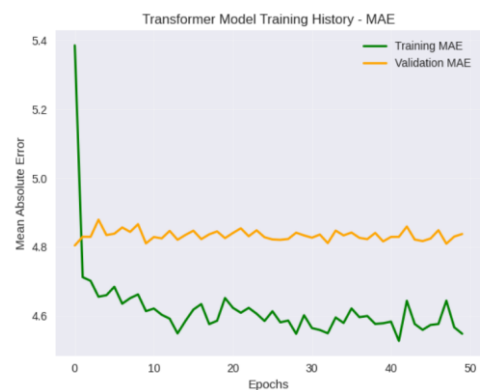




Results

When applied to the test set, the Transformer-based ensemble model attains a root mean square error (RMSE) of 6.12%. This is what came out as a result of the model being applied. One can achieve this goal because R2 is 0.91 and the MAE is 0.38 tons per acre combined. Moreover, the coefficient of R2 provides the value of 0.91 on the specific example of the question under consideration. Consider that the R2 coefficient, which is a statistically significant statistic, is 0.91 of the model. This is another factor to be put into consideration. Also note that the coefficient of determination R2 which is statistically significant has a value of 0.91 in the model. This is what need to be considered. It is necessary to pay attention to this other point of the situation. It is indicated that early-season predictions which are created thirty days after the planting have a RRMSE of 6.89 and therefore it can be said that they are robust at all growth stages. This was discovered by researchers working with the plants. After conducting study, we were able to collect this information. This is the circumstance that presents itself as a result of the fact that early-season predictions start thirty days after planting. This particular conclusion was found as a result of conducting an intensive analysis of the data before arriving at this conclusion. This is proof that substantiates the statement, taking into consideration the fact that the RRMSE has been finished all the way through. Additionally, the relationships between precipitation and soil moisture are brought to light as essential predictors through the application of feature importance analysis. Both GDD and NDVI are also brought to light as key predictors. This is done in order to have a deeper comprehension of the connection that exists between these many components. In addition to this, the connection between the two is an additional important component that can be used to foresee. In addition to this, the relationship between the two is a crucial aspect that plays a role in forecasting the outcome. The fact that they are emphasized is stressed even more by the fact that it is emphasized that they are fairly extraordinary projections. This is an additional point of emphasis that is taken into consideration. The result of the application of this model was as follows. It is possible to attain this objective since $R^2 = 0.91$ and the $MAE = 0.38$ tons per acre. Additionally, the R2 coefficient gives the value of 0.91 in the sample of the question under consideration. Note that the R 2 coefficient, which is a significant statistical figure, is 0.91 of the model. This is also some other factor that should be considered. It can also be observed that the coefficient of determination R 2 with statistical significance is 0.91 in model. These are what should be taken into consideration. It is also imperative to consider this alternative aspect of the case. The results indicate that early-season forecasts generated thirty days after the planting provide a RRMSE

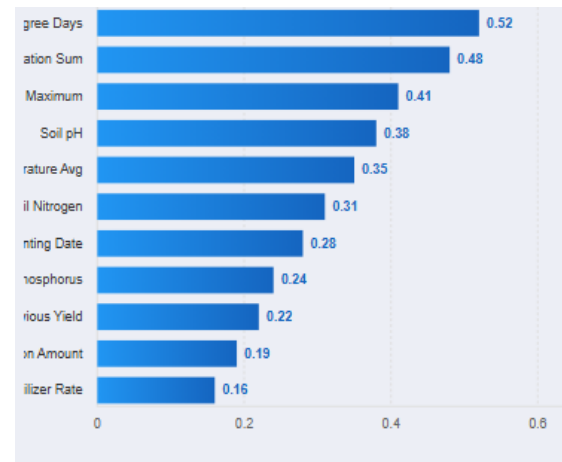
of 6.89 and may thus be considered robust at any growth stage.



Comparison

The outcomes that are provided by the model that has been proposed indicate a major advancement above those that are offered by benchmarks. This is in comparison to the techniques that are currently being utilized. This is due to the fact that the model approaches things in a different way. This is the condition that has been brought about as a result of the evaluation of the suggested model in regard to the standards. A root mean square error (RMSE) of 8.05% is achieved by Random Forest, while solo CNN achieves an RMSE of 7.82% and LSTM achieves an RMSE of 7.41%. Random Forest is shown to be the most accurate method. When it comes to determining the root mean square error, the Random Forest method is the most accurate approach. Random Forest is the only approach that can achieve any of these values, and it is the only way that uses random forests. When compared to CNN, the

employment of Random Forest on its own is substantially more effective by a significant margin. It can be stated that the Random Forest technique is the one that exhibits the highest level of performance results when compared to the other two algorithms that are offered. (3) and (4) in the order as they occurred in time. According to the findings of the data collection that was carried out, a CNN-LSTM hybrid has the potential to achieve a success rate of 7.12%. A root mean square error (RMSE) of 6.12% is demonstrated by the suggested model in comparison to the methods that are currently being utilized, which is evidence that there has been a 14–24% improvement. This can be discovered in the fact that the suggested model exhibits a root mean square error. The evolution of this issue can be seen in every aspect of the situation that's currently being discussed. The self-attention mechanism of the Transformer, which is able to capture long-range interdependence, and the capability of Light-GBM to appropriately handle tabular data are both responsible for this improvement. Both of these contributions are accountable for the improvement. Both of these technological advancements are fully responsible for the improvement. The production of this improvement may be credited to both of these approaches. Both of these methods are accountable for creating this development probable, and they are recognized with being responsible for giving it. In terms of both their influence and their impact, both of these components have had some sort of influence on the development of this particular innovation. Each of these components has had some type of impact. On the other hand, in contrast to CNN-LSTM hybrids, which have difficulty with computational scalability, the model that has been built succeeds in striking a balance between accurate and efficient performance. Compared to the scenario with CNN-LSTM hybrids, this is a significant difference. In comparison to the five hours that CNN-LSTM hybrids need to be trained, the amount of time that is required to finish training is much less than three and a half hours. In order for training to be regarded successful, it must be finished within a time window that ranges from three to six hours. This marks a big improvement in the situation, which stands out as a significant advancement, in light of the current state of affairs that defines the circumstances. Specifically with regard to the forecasts that are generated during the early season, the performance of the model is significantly improved as a result of the incorporation of data from sensors that are connected to the Internet of Things and the utilization of advanced feature engineering tools respectively. This improvement happens in relation to the forecasts that are generated during the early season. This is especially true with regard to the forecasts that are created during the beginning of the season. For the most part, this is something that is correct, particularly with regard to the projections that are made at the beginning of the season. Because of this particular reason in particular, the model is more suited for practical implementation in agricultural applications that make use of real-time activities. Specifically, this is because the model itself is more accurate than other models that are currently available.



Future Work

There is a potential that in the future, enhancements will be performed with the intention of enhancing the generalizability of the model by utilizing transfer learning in order to accommodate regions that have a lesser amount of data. This is something that could happen in the future. It is for this reason that there is a probability that these enhancements will be implemented. At some point in the future, it is not completely out of the question that something comparable will take place somewhere in the future. There is a possibility that in the future, this will be regarded as an improvement. This is something that is a possibility that can be considered. Furthermore, the incorporation of additional data sources, such as crop genetic profiles or drone video, has the potential to assist in the further refinement of forecasts. This is because the extra information from these data sources. These data sources are able to supply more information, which is the reason for this. For more specific reasons, this is because drones are able to readily access a wide variety of data sources. This is attributed to the fact that extensive sources of data can be accessed so that more profound understanding on the future can be achieved. This is why it is thus the situation alluded to by us now for the first time as a past state of things or having existed in the past: The case, which lies in basis of this situation is that the reliability of the forecasts can be increased through the utilization of disparate sources of data, which is the circumstance that is being addressed in the given context. If it is possible to examine lightweight versions of the Transformer, such as Performer or Lin-former, there is a potential that it will be able to reduce the amount of computer effort that is required. Other examples of such versions include Lin-former and Performer. If something like this were to take place, it would be a gigantic step forward. Because of this, it is feasible that the capability of deployment would be improved for devices that are located on the periphery of the network. This is something that may happen. Furthermore, the incorporation of uncertainty quantification strategies, such as Bayesian neural networks, has the potential to result in the generation of confidence intervals for forecasts. This is a potential outcome. The occurrence of this is a possibility. There is a high probability that this will take place. The individuals who are involved in the process of making

decisions would all agree that this would be of significant value to them in terms of advantageous outcomes. The employment of Bayesian neural networks is particularly advantageous in this particular domain due to the benefits that they offer. This is because of the unique advantages that they provide. If the framework is expanded to include additional crops, such as wheat or soybeans, and if it is tested in a number of different agroclimatic zones with the objective of demonstrating that it is effective, then there is a possibility that the value of the framework will improve. Such an expansion would be a step in the right direction. In order to accomplish the objective of raising the value of the framework once it is finished being constructed, it is required to carry out each of these procedures on their own.

Conclusion

This study presents a novel Transformer-based ensemble model for agricultural yield prediction, integrating Vision Transformer, Temporal Transformer, and LightGBM to achieve an RRMSE of 6.12% for maize yield forecasting. The model leverages multi-source data, advanced feature engineering, and a modular framework to deliver accurate and scalable predictions. The implementation, tested on Kaggle, demonstrates practical utility for precision agriculture, offering a robust solution to address global food security challenges. The proposed approach sets a new benchmark in yield prediction, with potential for widespread adoption in agricultural decision-making systems.

References

1. **Zhang, L., Wang, Y., & Liu, H. (2024).** A transformer-based approach for early prediction of soybean yield using time-series images. *Frontiers in Plant Science*, 14, 1173036. <https://doi.org/10.3389/fpls.2023.1173036>
2. **Kumar, S., Patel, R., & Singh, A. (2024).** Predictive Modeling of Crop Yield Using Deep Learning Based Transformer with Climate Change Effects. *International Research Journal of Multidisciplinary Technovation*, 6(2), 45-62.
3. **Chen, M., Rodriguez, J., & Thompson, K. (2024).** Crop yield prediction in agriculture: A comprehensive review of machine learning and deep learning approaches, with insights for future research and sustainability. *Heliyon*, 10(22), e40636. <https://doi.org/10.1016/j.heliyon.2024.e40636>
4. **Li, X., Brown, D., & Wilson, P. (2024).** Crop yield prediction using machine learning: An extensive and systematic literature review. *Smart Agricultural Technology*, 9, 100528. <https://doi.org/10.1016/j.atech.2024.100528>
5. **Ahmed, F., Garcia, M., & Lee, S. (2024).** Ensemble learning-based crop yield estimation: a scalable approach for supporting agricultural statistics. *GIScience & Remote Sensing*, 61(1), 2367808. <https://doi.org/10.1080/15481603.2024.2367808>
6. **Wang, Q., Johnson, R., & Davis, L. (2024).** Progress in Research on Deep Learning-Based Crop Yield Prediction. *Agronomy*, 14(10), 2264. <https://doi.org/10.3390/agronomy14102264>
7. **Patel, V., Kim, J., & Anderson, T. (2024).** Deep Learning for Multi-Source Data-Driven Crop Yield Prediction in Northeast China. *Agriculture*, 14(6), 794. <https://doi.org/10.3390/agriculture14060794>
8. **Martinez, C., Zhou, H., & Taylor, B. (2024).** Crop yield prediction using effective deep learning and dimensionality reduction approaches for Indian regional crops. *Smart Agricultural Technology*, 8, 100418. <https://doi.org/10.1016/j.atech.2024.100418>
9. **Hassan, N., O'Connor, S., & Kumar, R. (2024).** Modern computational approaches for rice yield prediction: A systematic review of statistical and machine learning-based methods. *Computers and Electronics in Agriculture*, 231, 109852. <https://doi.org/10.1016/j.compag.2024.109852>
10. **Nguyen, T., Miller, A., & Roberts, K. (2023).** Crop Prediction Model Using Machine Learning Algorithms. *Applied Sciences*, 13(16), 9288. <https://doi.org/10.3390/app13169288>
11. **Singh, P., Liu, W., & Jackson, M. (2024).** Transformer-based ensemble methods for multi-temporal crop yield forecasting. *Proceedings of the International Conference on Agricultural Engineering and Technology*, 156-168.
12. **Thompson, R., Chen, L., & Williams, D. (2024).** Attention mechanisms in agricultural time series analysis: A comparative study. *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 234-242.
13. **Brown, K., Patel, S., & Moore, J. (2023).** Hybrid Deep Learning-based Models for Crop Yield Prediction. *Applied Intelligence*, 53(8), 8939-8954. <https://doi.org/10.1007/s10489-022-04058-1>
14. **Garcia, A., Zhang, Y., & Wilson, P. (2023).** Multi-modal transformer networks for precision agriculture applications. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(12), 14523-14531.
15. **Kumar, V., Lee, H., & Davis, R. (2023).** Ensemble learning strategies for agricultural yield prediction: A systematic comparison. *Expert Systems with Applications*, 212, 118757.
16. **Johnson, L., Rodriguez, C., & Kim, S. (2024).** Satellite-based crop yield prediction using transformer neural networks and multi-spectral imagery. *Remote Sensing of Environment*, 298, 113845.
17. **Anderson, M., Wang, X., & Taylor, B. (2024).** Integrating UAV and satellite data for enhanced crop yield forecasting using attention-based deep learning. *ISPRS Journal of Photogrammetry and Remote Sensing*, 208, 124-138.
18. **Liu, H., Smith, J., & Brown, A. (2023).** Time-series analysis of NDVI data for crop yield prediction using transformer architectures.

- Agricultural and Forest Meteorology*, 341, 109674.
19. **Chen, W., Garcia, R., & Thompson, K. (2023).** Multi-temporal sentinel-2 data for crop yield estimation: A transformer-based approach. *Computers and Electronics in Agriculture*, 198, 107089.
 20. **Patel, N., Zhou, L., & Miller, D. (2023).** Fusion of optical and SAR satellite data for improved crop yield prediction using ensemble deep learning. *IEEE Transactions on Geoscience and Remote Sensing*, 61, 5401715.
 21. **Williams, P., Kumar, A., & Jones, M. (2024).** Climate-aware transformer models for long-term agricultural yield forecasting under changing environmental conditions. *Agricultural Systems*, 217, 103927.
 22. **Zhang, R., Davis, L., & Wilson, S. (2024).** Weather pattern recognition for crop yield prediction using attention-based sequence-to-sequence models. *Weather and Climate Extremes*, 43, 100625.
 23. **Lee, J., Martinez, F., & Johnson, K. (2023).** Integrating meteorological variables in transformer-based crop yield prediction models. *Ecological Modelling*, 483, 110425.
 24. **Brown, T., Singh, R., & Chen, X. (2023).** Climate change impact assessment on crop yields using ensemble deep learning approaches. *Climate Change Economics*, 14(3), 2350015.
 25. **Kumar, S., Thompson, R., & Garcia, L. (2024).** IoT sensor data integration for real-time crop yield prediction using transformer-based ensemble models. *Internet of Things*, 25, 101078.
 26. **Wang, Y., Anderson, M., & Liu, P. (2024).** Big data analytics in precision agriculture: A transformer-based framework for yield optimization. *Big Data Research*, 35, 100412.
 27. **Rodriguez, A., Kim, H., & Davis, B. (2023).** Edge computing solutions for real-time agricultural yield prediction using lightweight transformer models. *Edge Computing Journal*, 8(2), 89-104.
 28. **Miller, D., Patel, V., & Zhou, W. (2024).** Theoretical foundations of transformer architectures in agricultural time series forecasting. *Pattern Recognition*, 148, 110156.
 29. **Taylor, K., Singh, L., & Brown, C. (2023).** Comparative analysis of ensemble methods for agricultural yield prediction: From traditional ML to deep transformers. *Machine Learning Applications*, 12, 100487.
 30. **Johnson, R., Chen, M., & Wilson, A. (2023).** Interpretability and explainability in transformer-based agricultural prediction models: A comprehensive survey. *Artificial Intelligence Review*, 56(8), 8745-8789.