

Feasibility of machine learning-based rice yield prediction in India at the district level using climate reanalysis data

Djavan De Clercq, Adam Mahdi

University of Oxford

Abstract

Yield forecasting, the science of predicting agricultural productivity before the crop harvest occurs, helps a wide range of stakeholders make better decisions around agricultural planning. This study aims to investigate whether machine learning-based yield prediction models can capably predict Kharif season rice yields at the district level in India several months before the rice harvest takes place. The methodology involved training 19 machine learning models such as Cat-Boost, LightGBM, Orthogonal Matching Pursuit, and Extremely Randomized Trees on 20 years of climate, satellite, and rice yield data across 247 of India's rice-producing districts. In addition to model-building, a dynamic dashboard was built understand how the reliability of rice yield predictions varies across districts. The results of the proof-of-concept machine learning pipeline demonstrated that rice yields can be predicted with a reasonable degree of accuracy, with out-of-sample R², MAE, and MAPE performance of up to 0.82, 0.29, and 0.16 respectively. These results outperformed test set performance reported in related literature on rice yield modeling in other contexts and countries. In addition, SHAP value analysis was conducted to infer both the importance and directional impact of the climate and remote sensing variables included in the model. Important features driving rice yields included temperature, soil water volume, and leaf area index. In particular, higher temperatures in August correlate with increased rice yields, particularly when the leaf area index in August is also high. Building on the results, a proof-of-concept dashboard was developed to allow users to easily explore which districts may experience a rise or fall in yield relative to the previous year. The dashboard show that the model may perform better in some regions than in others. For instance, the absolute percentage error for predicted versus actual yields ranged from an average of 7.1% in districts in Uttarakhand to an average of 14.7% in Uttar Pradesh. This study underscores the potential for policymakers to consider scaling and operationalizing machine learning approaches to rice yield prediction in the context of agricultural early warning systems to deliver timely crop yield forecasts on a rolling basis throughout the season, thereby equipping agricultural decision-makers with the ability to make informed choices on irrigation scheduling, fertilizer application, and harvest planning to optimize crop output and resource use.

1 Introduction

1.1 The societal implications of accurate crop yield forecasting in India

Yield forecasting is the science of predicting agricultural productivity as measured by crop yield – the ratio of the total mass of the harvested product (such as rice) to the area used to cultivate the crop – before the harvest takes place, typically a few months in advance [1].

Pre-harvest prediction of crop yields is important in helping a wide range of stakeholders make better decisions around agricultural planning. For farmers, accurate crop yield forecasts can facilitate decision-making around what to grow and when to grow it [2]. In addition, near real-time monitoring of crop growth can inform the use of preventive measures such as irrigation and fertilization to boost agricultural productivity where needed [3]. For governments, yield prediction is relevant to the formulation of policies related to national food security, such as pricing policies for domestic markets, and policy decisions on the import and export of different crops [4].

Accurate crop yield forecasting may also enable better design of insurance products that mitigate climate risks and stabilize farmer incomes [5]. Weather-based crop insurance, for instance, uses a

weather index such as total precipitation to determine payments to farmers, meaning that insurance companies do not need to visit farmers to assess damages and arbitrate claims. Rather, if the weather reaches a certain threshold, rapid automatic payments can be distributed to farmers, who avoid the need to sell assets to survive due to adverse climate events [6].

The need for accurate information on crop yields is particularly important in countries like India, where the agricultural sector provides livelihoods for hundreds of millions of farmers, with 70% of rural households depending on agriculture for their main source of income. One of India's major staple crops is rice, which contributes to 30% of calories consumed in India and is a key export commodity for the country [7]. India cultivates rice on about 45 million hectares of land, with a total production of 178 million metric tonnes in 2020 [8]. In addition, the distribution of monsoon rainfall, which is a major source of water for rice cultivation, has become erratic in recent years due to climate variability [9, 10]. In such contexts, crop yield predictions may be able to supplement agricultural early warning systems (AEWSs) that give advanced notice of potential risks to crop productivity, enabling preemptive action in affected areas. Previous research has shown that current agricultural monitoring systems lack robust crop yield and crop production forecasts, as well as the operationalisation of such methods at scale [11].

1.2 Overview of approaches and variables used to model crop yields

Crop yield prediction is a challenging problem in precision agriculture, as final yields depend on a variety of factors such as weather, climate, soil, seed type, and agronomic practices such as irrigation and fertilizer use [12].

This complexity is evident from the variety of variables included and methods applied in the growing body of literature on crop yield forecasting. For example, recent examples in literature involving deep learning approaches include corn and soybean yield forecasting in the US based on convolutional neural networks (CNNs) and recurrent neural networks (RNNs) [13], soybean yield forecasting in Argentina based on deep transfer learning [14], and vineyard grape yield estimations based on CNNs [15]. Recent examples based on machine learning approaches include sugarcane yield prediction using random forests [16], prediction of wheat, barley, and canola yields in Western Australia using random forest [17], yield forecasting of spring maize in Pakistan based on LASSO regression and support vector machine [18], and Jojoba yield prediction in Israel based on gradient boosted regression trees [19]. Other examples of the machine learning approaches that have been applied to yield prediction have been summarized in a systematic literature review, which also analyzed the variables most frequently included in crop yield prediction studies. Across 50 studies between 2008 and 2019, features used as predictors of yield have included temperature, soil type, rainfall, humidity, pH-value, NDVI, wind speed, and more [2].

For studies specific to rice, the staple crop of over half the world's population [20], a number of approaches have been applied to yield forecasting in recent years. Recent examples include rice yield prediction for 81 counties in southern China based on recurrent neural networks [21]; application of the ecological distance algorithm to model rice yields [22]; field and county-level rice yield prediction based on synthetic aperture radar (SAR), optical and meteorological data [23]; random forest yield prediction based on high-resolution imagery collected from unmanned aerial vehicles (UAVs) [24]; simulation of yields using the Cropping System Model-CERES-RICE [25]; pixel-scale rice yield prediction in South Korea based on a combination of deep learning and crop models [26]; rice yield estimation at 500m spatial resolution based on gradient boosted regression and vegetation indices derived from the Moderate Resolution Imaging Spectroradiometer (MODIS) [27]; and rice paddy yield prediction using sentinel-based optical and SAR data in India based on random forest [28]. To the authors' knowledge, there has been less research on rice prediction at the district-level in India.

1.3 Research contributions

This study marks an advancement in the field of rice yield prediction in India by building a proof-of-concept approach capable of predicting rice yields at the district level for 247 rice-producing districts across India.

First, a novel combination of data sources is used to predict Indian rice yields. These include data from ERA5, a climate re-analysis product developed by the European Centre for Medium Range Weather Forecasts (ECMWF), which combines observations with modelled data to provide hourly

data on atmospheric, land-surface, and sea-state parameters globally [29]. Vegetation data was derived from the MODIS sensors on-board NASA’s TERRA and AQUA satellites, and a cropland mask (CROPGREIDS) was used to filter earth observation data according to where rice is grown at a pixel level. Collectively, these data enhance the model’s ability to capture the intricate effects of climate variables and vegetation health indices on crop yields.

Secondly, the study creates a spatially matched dataset for rice crop yields in India, tackling a common problem in yield prediction research: a lack of training datasets where yield data is accurately aligned with specific geographic locations. Typically, discrepancies between the names and boundaries in geographic datasets (like shapefiles) and official agriculture statistics can hinder the effective use of data in research. Fuzzy matching algorithms were used to align Indian shapefiles (which may have variations in district names) with official rice yield data from the Indian government to produce a dataset where yield information is precisely matched to its geographic location. This matching is important for researchers aiming to spatially aggregate earth observation data in a manner consistent with available yield data. Making this matched dataset open source specifically aids rice-related research in India, enabling more accurate earth observation studies by providing a reliable foundation for correlating satellite data with actual agricultural outputs.

Thirdly, the study provides a new benchmark for district-level rice yield prediction in India based on predictions made exploring the predictive effectiveness of 19 machine learning models such as LightGBM, Bayesian Ridge Regression, and others. This methodological approach facilitates the identification of the most effective models for rice yield prediction at the district level. Previous literature on crop yield prediction in India has largely focused on using a narrower range of algorithms such as random forest or support vector machine [2]. A detailed evaluation of model performance is provided, coupled with the interpretability provided by SHAP values.

Fourth, beyond model development, an interactive dashboard tool was developed, not only to visualize the yield predictions across each of India’s districts, but also to allow for detailed diagnosis of model performance. It serves as a practical tool for model evaluation, offering insights into regional performance variations and facilitating the identification of areas where predictions may be improved.

Overall, this study offers evidence that scalable crop yield prediction models have the potential to be integrated into agricultural early warning systems in India, which, as noted in previous research, currently lack such forecasting capabilities and the means to operationalize these methods at a scale that contributes towards more resilient agricultural practices and food security planning [11].

2 Methods

2.1 Study region and brief overview

In this study, climate and remote sensing data were used as predictors to model rice yields for the kharif season (wet summer monsoon season) from 2001 to 2020 at the district level in India (India consists of 36 states and 684 districts). In India, more than half of the annual rice crop is grown during kharif [30], a season which is characterized by high temperature, high humidity, and medium to high rainfall [31]. Kharif season rice is typically sowed between the start of June to the end of August and harvested between the end of September to early January, depending on the region. During the 2019-2020 season, harvesting of Kharif rice was completed in February 2020.

The methods applied in this study can be summarized as follows. Firstly, 20 years of data on climate, vegetation, and rice yields were ingested programmatically via API from various sources. Second, climate and vegetation data were pre-processed and aggregated to the same level as historical rice yields prior to development of machine learning models. Third, a range of models including Bayesian ridge regression and LightGBM were trained, evaluated, and interpreted. Fourth, an interactive dashboard tool was developed to visualize the results of the algorithms to provide a spatially explicit view of model results and potential modelling errors. Lastly, the implications of the modelling results were discussed with regards to their inclusion in agricultural early warning systems.

The methods outlined in this research are fully reproducible. All data, python code, and dashboards can be found on GitHub.

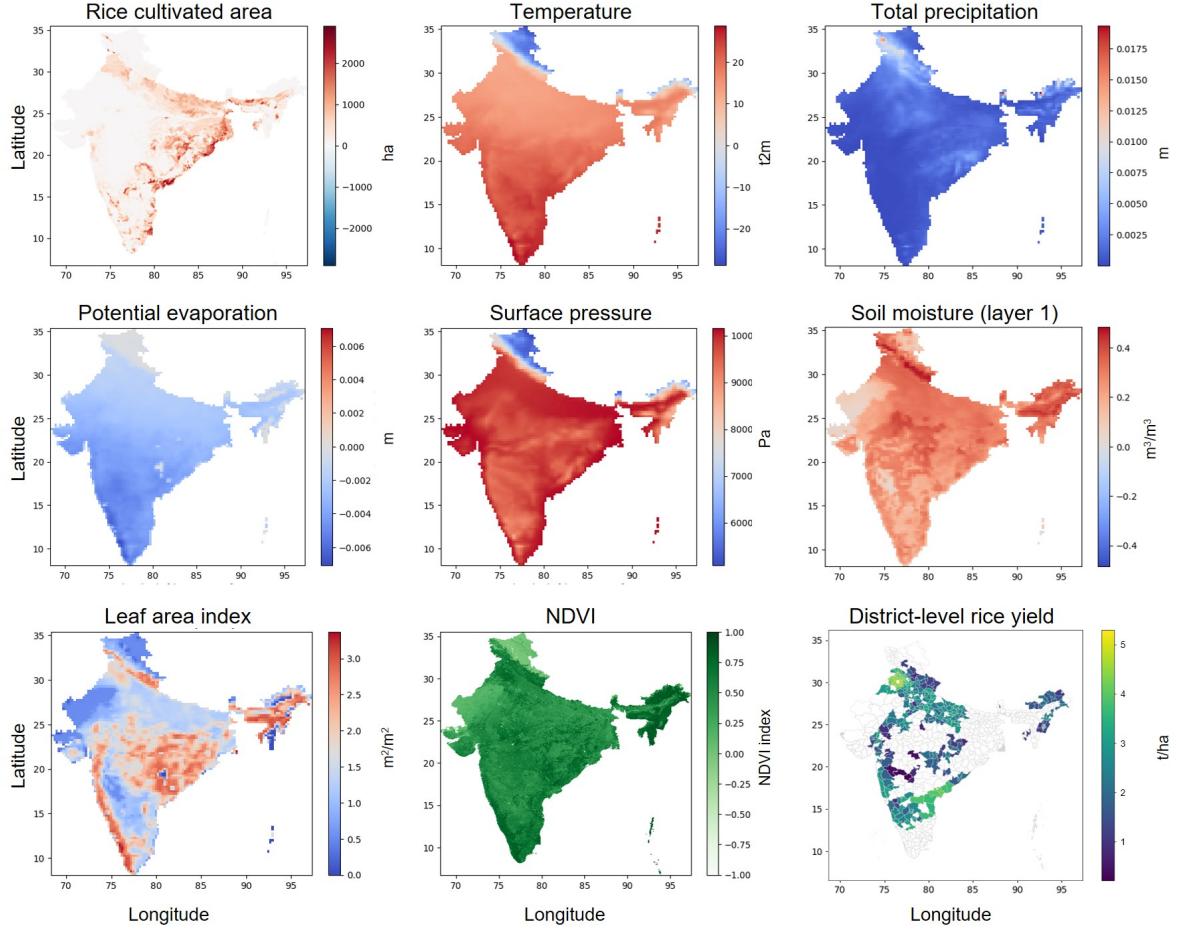


Figure 1: Overview of the geospatial data used in this research. The panels show, respectively: cultivated rice area in India; snapshots of temperature, total precipitation, potential evaporation, surface pressure, soil moisture, and leaf area index from the ECMWF (average values of January 2022); NDVI data from NASA’s MODIS (average values of January 2022); and district-level rice yield data from India’s Ministry of Agriculture and Farmers Welfare (in 2020).

2.2 Data collection

As shown in Table 1, a range of datasets were used in the rice yield modelling, including climate reanalysis data, remote sensing data, district-wise historical rice yield data, and cropland masks. A visual overview of the geospatial data used in this research are also provided in Figure 1.

Climate reanalysis data. Daily climate reanalysis data on temperature, potential evaporation, surface pressure, leaf area index, total precipitation, and soil water content was obtained from ERA5 data from the European Center for Medium-Range Weather Forecasts (ECMWF), which provides global estimates of surface and atmospheric parameters since 1950 at a resolution of approximately 30°*30 km [32]. Climate reanalysis data, which are often freely available, provide temporally and spatially homogenous data [33], which makes them suitable for applications such as crop yield prediction in contexts where in-situ weather station measurements are inadequate or incomplete. In addition, weather stations vary in their accuracy and generally record a limited number of variables, such as rainfall, temperature, pressure, and wind speed; variables that are more technically demanding to measure, such as humidity and solar radiation, may be lacking [34].

The climate variables used in this study were selected due to their influence on rice yields. An extensive body of research has shown that rice growth is affected by factors such as soil water content [35], temperature [36], potential evaporation (as a proxy for transpirational demand) [37], surface pressure [38], and precipitation [39]. These variables can impact rice development across phenological

stages. Past case studies in India have shown the approximate number of days taken for each stage: the sowing to tillering phase (P1) can range from 30 to 60 days, and the rate of tillering tends to increase under higher temperatures; the tillering to panicle initiation phase (P2) can range between 42 to 49 days; the panicle initiation to flowering phase (P3) can range from 12 to 28 days; the flowering to milk phase (P4) can range from 7 to 20 days; and the mil to physiological maturity phase (P5) can range from 17 to 31 days [40]. One study showed reported a linear relationship between the days taken from sowing to flowering and average air temperature [41].

Remote sensing data. Normalized Difference Vegetation Index (NDVI) is a dimensionless index that describes the difference between visible and near-infrared reflectance of vegetation cover, and can be used to estimate the density of green on an area of land [42]. To determine the density of green on a patch of land, the wavelengths of visible and near-infrared sunlight reflected by the plants are observed. NDVI values range from -1 to +1; higher values of NDVI imply healthy and dense vegetation, whereas lower NDVI values indicate sparser vegetation. NDVI data was obtained from the Moderate Resolution Imaging Spectroradiometer (MODIS) onboard NASA's Terra and Aqua satellites, due to their wide coverage and temporal resolution. There are several examples in academic literature of using NDVI to investigate the progress of crops, such as for wheat in Argentina [43], cereals in Europe [44], and rice in Vietnam [45]. NDVI data was masked using CROPGRIDS, a global, geo-referenced dataset providing information on areas for 173 crops circa the year 2020, at a resolution of 0.05° (~ 5.55 km at the equator) [46].

Yield data. District-level rice production and yield data from 1995 to 2021 for 367 districts were obtained from the APY dataset of the Directorate of Economics and Statistics in India's Ministry of Agriculture and Farmers Welfare [47]. In this dataset, the year denotes the year in which the crop was harvested. For kharif season rice, the sowing is in the previous calendar year [48]. The raw yield data were aligned with a shapefile representing Indian districts. This alignment was achieved using the FuzzyWuzzy library in Python, which employs the Levenshtein distance to address the challenge of slight spelling discrepancies in the district names between the yield data and the shapefile [49].

Data used in this analysis were programmatically downloaded via API and automated Python scripts. ERA5 data was ingested via the CDS API, while NDVI data was ingested via USGS' AppEEARS API.

2.3 Data pre-processing

Climate data from ERA5 in NetCDF format over a bounded area comprising India was clipped to the Indian country boundary. NDVI data from NASA's AQUA MODIS satellite in NetCDF format were clipped to the Indian country boundary and masked with a rice cropland layer.

Next, the climate variables and NDVI were aggregated to the district level based on zonal statistics. The vector geometry data for India's ADM2 (district-level) boundaries which raster pixels were aggregated to were obtained from the Database of Global Administrative Areas (GADM) [50]. District-level yield data from APY was then merged to the climate and remote sensing data aggregated at the district level to produce a spatially consistent geodataframe. Yield outliers beyond three standard deviations were removed as they were assumed not achievable at the district level in India [27].

Feature engineering was conducted to produce monthly averages for the climate and NDVI parameters for every month between May and November, corresponding to the full sowing and growing period for kharif rice [27]. This process was repeated for all variables to produce a set of 52 features used as input for the modelling. The months selected for climate and NDVI feature aggregation were chosen to reflect the full range of rice growth stages, including the grain filling, vegetative, and reproductive stages [51].

2.4 Model development and interpretation

This study developed and tested the performance of multiple rice yield prediction models based on a variety of machine learning models. These included LightGBM [52], an efficient and distributed gradient boosting framework that uses tree-based learning, Bayesian ridge regression [53, 54, 55], which has been recognized for its ability to deal with hierarchical data structures [56], gradient boosting regression [57], random forest [58], Huber regression [59], decision tree regression [60], elastic net regression [61], AdaBoost [62], orthogonal matching pursuit [63], and extremely randomized trees [64].

Table 1: A range of agronomically-relevant datasets were used as predictors of the target variable (district-level rice yield in India); rice areas masks were used to filter NDVI data by rice-growing area

Data type	Parameter	Description	Unit	Source
Climate reanalysis	Potential evaporation (pev)	A measure of the extent to which near-surface atmospheric conditions are conducive to the process of evaporation.	m	ECMFW (ERA5)
	2m-temperature	The temperature of air at 2m above the surface of land, sea or inland waters. 2m temperature is calculated by interpolating between the lowest model level and the Earth's surface, taking account of the atmospheric conditions.	K	ECMFW (ERA5)
	Total precipitation	The accumulated liquid and frozen water, comprising rain and snow, that falls to the Earth's surface. It is the sum of large-scale precipitation and convective precipitation.	m	ECMFW (ERA5)
	Leaf area index, low vegetation (LAI)	The surface area of one side of all the leaves found over an area of land for vegetation classified as 'low'. 'Low vegetation' consists of crops and mixed farming, irrigated crops, short grass, and more.	$m^2 \cdot m^{-2}$	ECMFW (ERA5)
Climate reanalysis	Total precipitation	The accumulated liquid and frozen water, comprising rain and snow, that falls to the Earth's surface. It is the sum of large-scale precipitation and convective precipitation.	m	ECMFW (ERA5)
	Volumetric soil water (SWVL1)	The volume of water in soil layer 1 (0 - 7cm, the surface is at 0cm)	$m^3 \cdot m^{-3}$	ECMFW (ERA5)
Remote Sensing	Normalized Difference Vegetation Index (NDVI)	A dimensionless index that describes the difference between visible and near-infrared reflectance of vegetation cover, and can be used to estimate the density of green on an area of land	(-)	NASA EOSDIS (AQUA MODIS)
Crop mask	Rice crop mask	A comprehensive global, geo-referenced dataset providing information on areas for 173 crops circa the year 2020, at a resolution of 0.05° (5.55 km at the equator).	ha	CROP

The models above were trained on district-level data for 2001 to 2018 (4,606 observations), and validated on out-of-sample test data for 2019 and 2020 (502 observations). The data was split in a manner that reflects how yield prediction models may be used in practice, avoiding random splits in favor of chronological splits to help ensure the model's robustness to future, unseen data. This out-of-sample approach to testing regression models with temporal dependency has been shown to be more robust than cross-validation approaches tailored to time series problems [65].

The top-performing models were evaluated based on three out-of-sample performance measures including R2, Mean Absolute Error (MAE), and Root Mean Square Error (RMSE). Also reported were Mean Squared Error (MSE), Mean Absolute Percentage Error (MAPE), and Root Mean Squared Logarithmic Error (RMSLE). In addition, model results were evaluated based on prediction error plots, residual plots, and spatio-temporal plots of prediction error to evaluate potential model bias (for example, better model performance for certain rice-growing regions of India). Lastly, Shapley Additive exPlanations (SHAP) were used to explore the impact of features on model output [66]. SHAP values are a model-independent methodology used for quantifying the significance of features in predictive modeling. The SHAP of feature j for observation $\phi_j(\mathbf{x})$ is defined as

$$\phi_j(\mathbf{x}) = \sum_{S \subseteq \{1, \dots, M\} \setminus \{j\}} \frac{|S|! \cdot (M - |S| - 1)!}{M!} (f(\mathbf{x}_{S \cup \{j\}}) - f(\mathbf{x}_S)) \quad (1)$$

where j is the feature evaluated, M the total number of features, S a subset of the full feature set $\{1, \dots, M\}$ that does not include the feature j , \mathbf{x}_S a subset of features in S , and f the model's prediction function [67].

2.5 Computation

Data ingestion, pre-processing, and modelling was conducted in a conda-based python environment with a diverse set of python libraries. Data processing and geospatial operations were carried out using python libraries including numpy, xarray, pandas, rasterio, rasterstat, and geopandas. Modelling and visualization was conducted using python libraries including scikit-learn, pycaret, matplotlib, and seaborn.

2.6 Interactive visualization for model evaluation and decision-making

Visual dashboards can serve as a helpful tool for making the outputs of analytical models more comprehensible to stakeholders by converting predictive analytics into easily interpretable visual formats. This transformation is particularly beneficial for users from various backgrounds, including farmers, policymakers, and researchers, as it facilitates their engagement with and comprehension of the data. For example, platforms like Streamlit and PowerBI have been used to build visual web applications and dashboards across diverse fields, such as bioinformatics, bacterial testing, financial auditing, twitter sentiment analysis, credit card fraud detection, drug target prioritization, and pharmaceutical sales forecasting [68, 69, 70, 71, 72, 73].

In this study, we utilize the dashboard to not only visualize the model's predictions, such as forecasted yields in each Indian district, but also to offer insights into the model's diagnostic aspects, like its varying predictive accuracy across different regions. We developed two distinct visual dashboards: the first provides a clear, spatially detailed overview of the model's forecasts, and the second aids in identifying and understanding any potential errors in the model's performance.

3 Results

This study assessed the feasibility of predicting pre-harvest rice yields in India using machine learning models, with a focus on satellite and climate reanalysis data. The results, which are detailed in the subsequent sections and compared to other global studies of rice yield prediction, show that the models achieve strong performance across several metrics in out-of-sample tests. Such results affirm the potential of these models for rice yield forecasting and provide a benchmark for predictive precision in the field.

3.1 Overview of out-of-sample model performance

Table 2 summarizes the out-of-sample (validation set) performance across the models tested: R2, MAE, and MAPE values of up to 0.82, 0.29, and 0.16 respectively were achieved. Compared to out-of-sample results reported in previous literature on rice yield prediction in different parts of the world, the models perform well. For instance, one study which developed rice yield prediction models for China based on support vector machine regression, neural networks, and random forest, achieved R2 values ranging from 0.24 to 0.31 and MAE values ranging from 0.58 to 0/66 t/ha [74]. Another study estimating rice yields in Vietnam’s Mekong Delta reported out-of-sample MAE values ranging from 0.46 to 0.55 t/ha for Winter and Summer rice models [75]. A study on county-level rice yield prediction in China’s Jiangsu province reported out-of-sample R2 values of 0.39 to 0.59 on an independent holdout set [23]. A study on pixel-scale rice yield prediction in South Korea reported test-set R2 values of 0.80 [26]. One image-driven yield prediction study reported test-set R2 values of 0.65 [76]. A study using multi-temporal UAV-based multispectral vegetation indices reported test set R2 values of up to 0.80 [77].

Table 2: Model performance on out-of-sample test data shows that the top three (based on R2 and MAPE) models include Random Forest, CatBoost, and Light Gradient Boosting. Experiment 1 (“all features”) shows results for model runs including climate/satellite-derived features and the “Year” and “District” features. Experiment 2 shows results for model runs trained exclusively using climate and satellite data observations (“EO features only”).

Model	Experiment 1 – all features				Experiment 2 – EO features only			
	MAE	RMSE	R2	MAPE	MAE	RMSE	R2	MAPE
Random Forest Regressor	0.31	0.41	0.80	0.16	0.45	0.56	0.63	0.25
CatBoost Regressor	0.29	0.39	0.82	0.18	0.43	0.53	0.67	0.24
Light Gradient Boosting Machine	0.31	0.41	0.80	0.19	0.44	0.56	0.64	0.25
Extreme Gradient Boosting	0.33	0.43	0.78	0.19	0.44	0.58	0.61	0.23
Orthogonal Matching Pursuit	0.33	0.46	0.76	0.20	0.76	0.96	-0.08	0.49
Decision Tree Regressor	0.41	0.56	0.63	0.20	0.58	0.81	0.24	0.33
Bayesian Ridge	0.33	0.46	0.76	0.21	7.51	11.81	-161.70	4.15
Gradient Boosting Regressor	0.32	0.41	0.80	0.21	0.50	0.62	0.55	0.28
Ridge Regression	0.34	0.47	0.75	0.21	0.61	0.78	0.30	0.37
Huber Regressor	0.33	0.46	0.75	0.21	0.75	0.95	-0.06	0.49
K Neighbors Regressor	0.39	0.51	0.70	0.21	0.55	0.70	0.44	0.31
Linear Regression	0.36	0.48	0.73	0.21	0.61	0.78	0.29	0.36
AdaBoost Regressor	0.45	0.55	0.65	0.27	0.63	0.75	0.34	0.40
Passive Aggressive Regressor	0.48	0.62	0.56	0.31	0.70	0.91	0.04	0.51
Elastic Net	0.66	0.81	0.24	0.41	0.74	0.94	-0.04	0.49
Lasso Regression	0.80	0.99	-0.13	0.53	0.74	0.94	-0.03	0.49
Lasso Least Angle Regression	0.80	0.99	-0.13	0.53	0.76	0.96	-0.07	0.48
Dummy Regressor	0.80	0.99	-0.13	0.53	0.80	0.99	-0.13	0.53
Least Angle Regression	1.90	2.35	-5.42	0.99	2.33	2.98	-9.34	1.21

The results also perform well compared to studies which only reported in-sample performance metrics. One study on rice yield modelling in Bangladesh reported in-sample R2 values ranging from 0.44 to 0.91; out-of-sample performance was not reported [78]. Another study on rice yield in China report in-sample R2 values of 0.77, lower than the out-of-sample R2 performance achieved in this study of 0.82 [79]. For rice yield prediction in the Philippines, one study reported an in-sample RMSE of 0.46 t/ha [80]. Another study using drones reported in-sample R2 values of 0.60 to 0.81 for rice yield prediction in Japan based on NDVI [81].

3.2 Errors, residuals, and SHAP value analysis

As shown in Figure 2, observed and simulated yields show a high level of agreement for some of the top-performing models including the random forest, CatBoost, and LightGBM regressors. In addition, the residual plot residual plots show that the majority of both training set and test set observations are randomly dispersed along the horizontal axis, indicating a reasonable low level of bias and homoscedasticity. The distribution of residuals here is roughly centred around zero but with some skewness, indicating the potential presence of outliers.

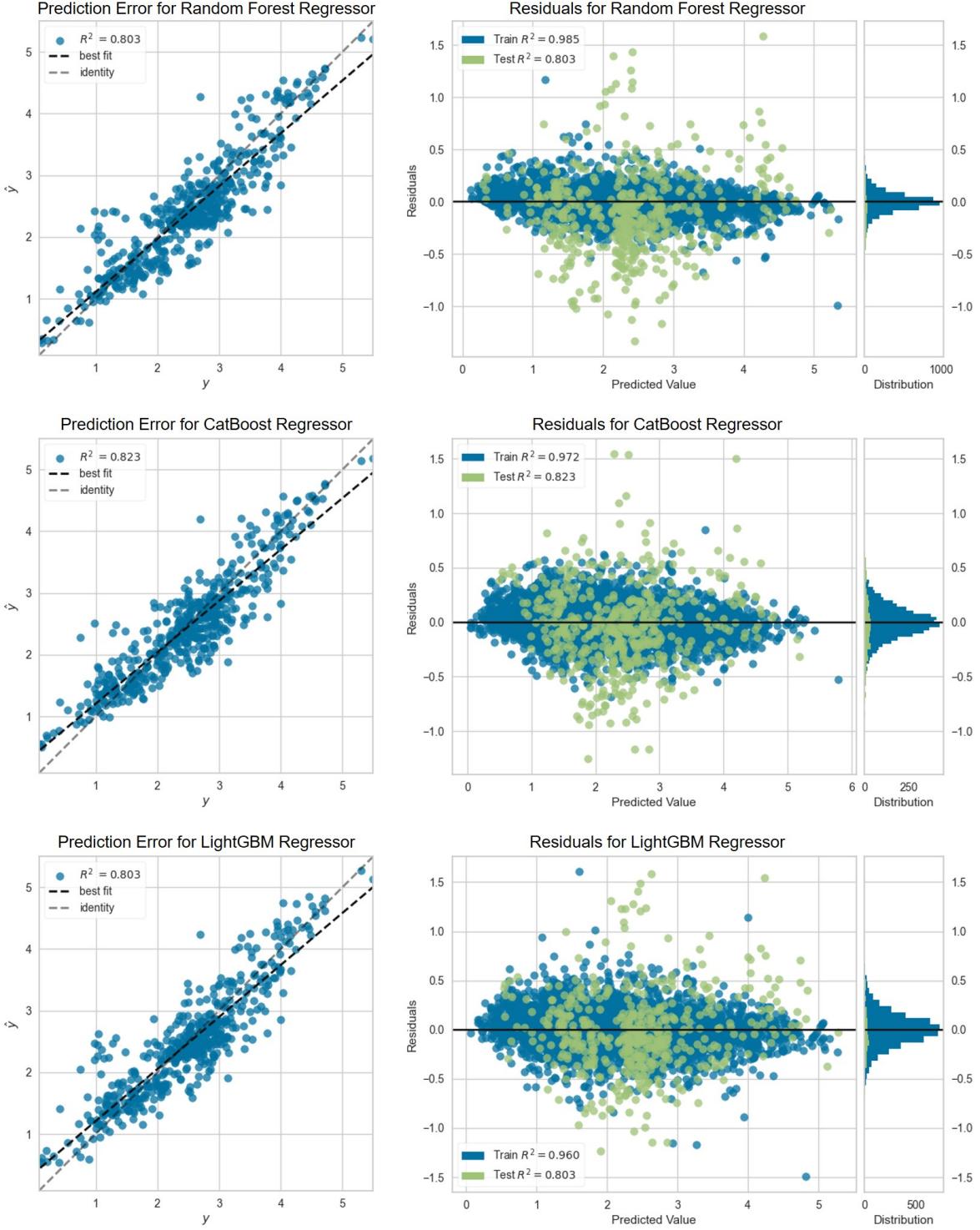


Figure 2: Comparative evaluation of selected models, including Random Forest, CatBoost, and LightGBM regressors using prediction error and residuals analysis. Two years of observations (502 observations in total) were used for the out-of-sample validation data, on which the Random Forest, CatBoost, and LightGBM models have test R^2 values of 0.80, 0.82, and 0.80 respectively. Residuals are mostly centered around zero, but CatBoost shows a skewness in error distribution. The histogram of residuals indicates Random Forest and CatBoost have a tighter error distribution compared to LightGBM's broader range.

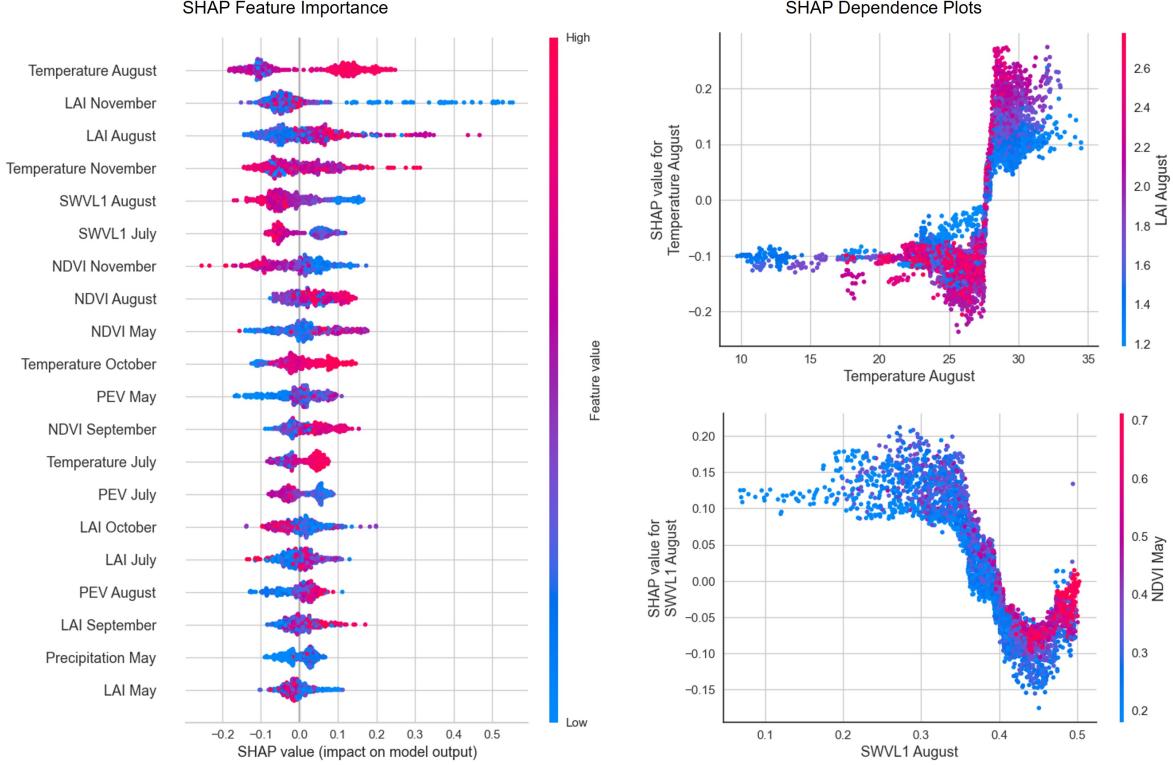


Figure 3: Interpretation of SHAP values for selected features in the rice yield prediction random forest model. The SHAP feature importance plot (left panel) exhibits the impact of various features on the model’s output. Higher SHAP values indicate a greater influence on the predicted yield. The coloring on the feature importance plot represents the value of the feature for each data point. Blue points indicate low feature values, while pink points represent high feature values. This color gradient allows us to visualize not only the impact (magnitude of the SHAP value) each feature has on the model output but also the distribution of the feature’s values. For instance, when examining ‘Temperature August’, we can see a mix of pink and blue points across a range of SHAP values, indicating a diverse range of temperatures in August within the dataset and how these varying temperatures correlate with the rice yield prediction. The top right panel presents a SHAP dependence plot for temperature in August, illustrating a correlation between higher temperatures and increased SHAP values for rice yield. The intensity of the color indicates the interaction effect, with a notable interaction with LAI in August, as higher LAI values (depicted in red) intensify the impact of temperature on yield. The bottom right panel depicts a SHAP dependence plot for soil water volume (“SWVL1”) in August, showing the relationship between SWVL1 values and SHAP values. This plot reveals that certain values of SWVL1 are associated with lower or higher SHAP values, indicating its varying influence on yield predictions, with the color intensity representing the interaction with NDVI in May.

In addition to the errors and residuals, the SHAP summary plots in Figure 3 concisely display the magnitude, prevalence and direction of a variable’s effect on final rice yield. The plots reveal that important variables include temperature, soil water volume (“SWVL1”), NDVI, and LAI in selected months. The importance ranking of these variables corroborates previous findings that that factors such as soil water content, temperature, and NDVI are important factors in estimating rice growth [35, 36].

A closer analysis of the specific impacts of features on rice yield is shown in the right-hand side panels of Figure 3. For instance, increases in temperature in August, which coincides with the sowing to panicle initiation phases of rice growth, are associated with higher yields, corroborating previous findings in India that above average yields may be associated with higher maximum temperatures [40]. The full set of SHAP plots for all features is available in the analysis output on GitHub.

3.3 Interactive visualization tool

In addition to the SHAP plots, which provide insight into how variables drive yield outcomes, two visual dashboards were developed to (a) provide an easy-to-understand, spatially explicit summary of model predictions, and (b) to help to identify potential biases in model performance. These are shown in Figure 4 and Figure 5.

For example, in Figure 4 for the year 2020 (one of the test set years) the dashboard shows that districts in the state of Chhattisgarh with expected yield increases relative to the previous year included Jashpur, Korba, and Koriya, where yields were expected to increase by 52%, 22%, and 22%, respectively. In another example, districts in the state of Gujarat such as Kheda and Sabar Kantha were expected to see yields decrease by 27% and 21% respectively according the model.

In addition to a visual representation of the predictions, the dashboard also provides a spatial view of prediction errors in order to more easily identify areas where model predictions may be inaccurate (Figure 5). For instance, the dashboard shows that on average, the absolute percentage error for predicted versus actual yields in ranged from an average of 7.1% in districts in Uttarakhand to an average of 14.7% in Uttar Pradesh, implying that the model may perform better in some regions than others.

3.4 Limitations and future research directions

Below, we present a non-exhaustive list of potential future research directions to build on the results of this research, including: refining model accuracy by exploring additional variables; communicating the outputs of yield prediction models in early warning systems by leveraging large language models; modeling how policymakers can help to disseminate the yield predictions-based early warning tools; and combining yield prediction modeling results with optimization approaches to support to anticipatory aid allocation efforts.

Modeling enhancements to increase predictive power. Increased predictive power of rice yield model may potentially be achieved by incorporating a wider array of agronomically relevant climatological variables. Prior investigations have highlighted the influence of thermal extremes, precipitation extremes, daytime humidity variations, and solar radiation on rice yields [40]. In addition, agricultural yield models – particularly those that are grounded in time series analysis – may stand to gain from incorporating additional auto-regressive elements, rolling averages, or cumulative indices, such as total sunshine hours or cumulative rainfall since rice sowing [82]. Beyond remotely sensed variables, yield models may also benefit from the inclusion of variables related to socio-economic changes in farmer populations, cultural practices that affect rice cultivation, and market or policy-related shifts that provide incentives for farmers to cultivate their crop in different ways [83, 84].

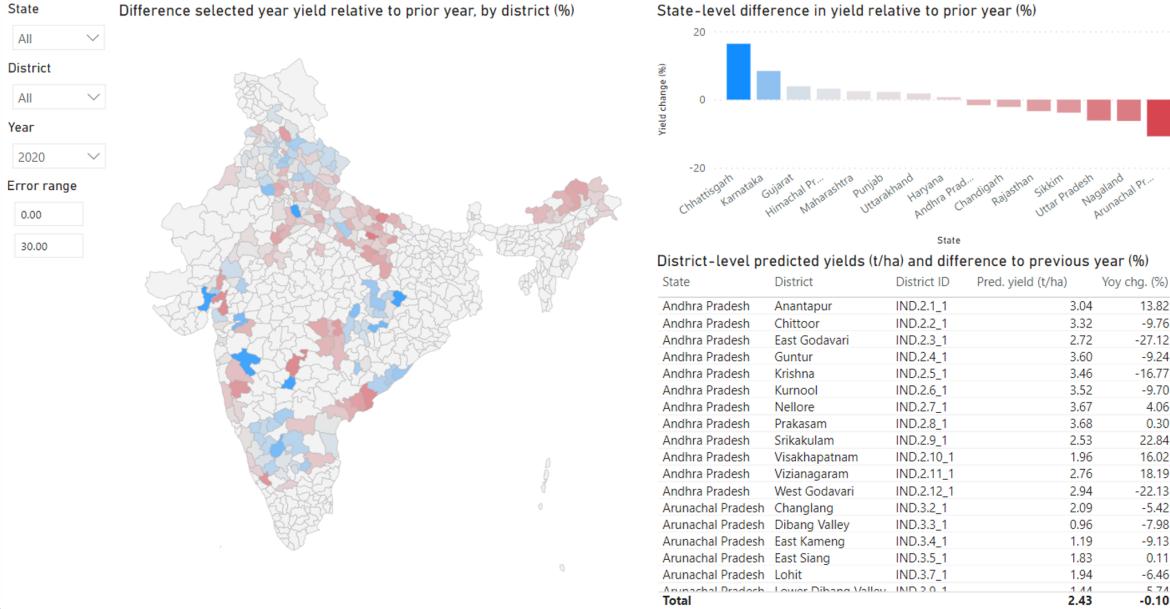
However, the augmentation of the feature set should also be approached with caution, as an expanded variable set can inadvertently lead to over-fitting, introduce feature redundancy, and complicate model interpretability. In scenarios where model clarity and understandability are important, a more conservative approach with a reduced set of features might be advisable to streamline the interpretive process, allowing for clearer insights and more straightforward decision-making [60]. Approaches which balance model interpretability with overall accuracy may be more helpful to decision makers.

Leveraging large language models to communicate yield prediction results. Additional research is needed on how access to climate and yield anomaly information via user-tailored interfaces can help mitigate meteorological shocks for agricultural communities. Studies have shown that inadequate access to climate information in South Asia has been observed as a factor driving perceived losses in farming communities [85].

Large language models (LLMs) may have potential to enhance the utility and accessibility of crop yield prediction models. By integrating LLMs with rice crop yield models, the data these yield models generate can be transformed into concise reports disseminated through channels such as text message-based farmer advisories or agricultural extension services. LLMs are able to process both structured and unstructured data, such as summarizing tabular data on yield predictions for various regions and translating these summaries into different local languages, thus increasing their accessibility and ease of understanding [86, 87].

Moreover, LLMs may have the potential to offer personalized agronomic advice based on location-specific yield predictions. Their proven capability in handling expert-level tasks across various fields,

PREDICTED RICE YIELDS IN RICE-GROWING REGIONS OF INDIA RELATIVE TO PREVIOUS YEAR



PREDICTED RICE YIELDS IN RICE-GROWING REGIONS OF INDIA RELATIVE TO PREVIOUS YEAR

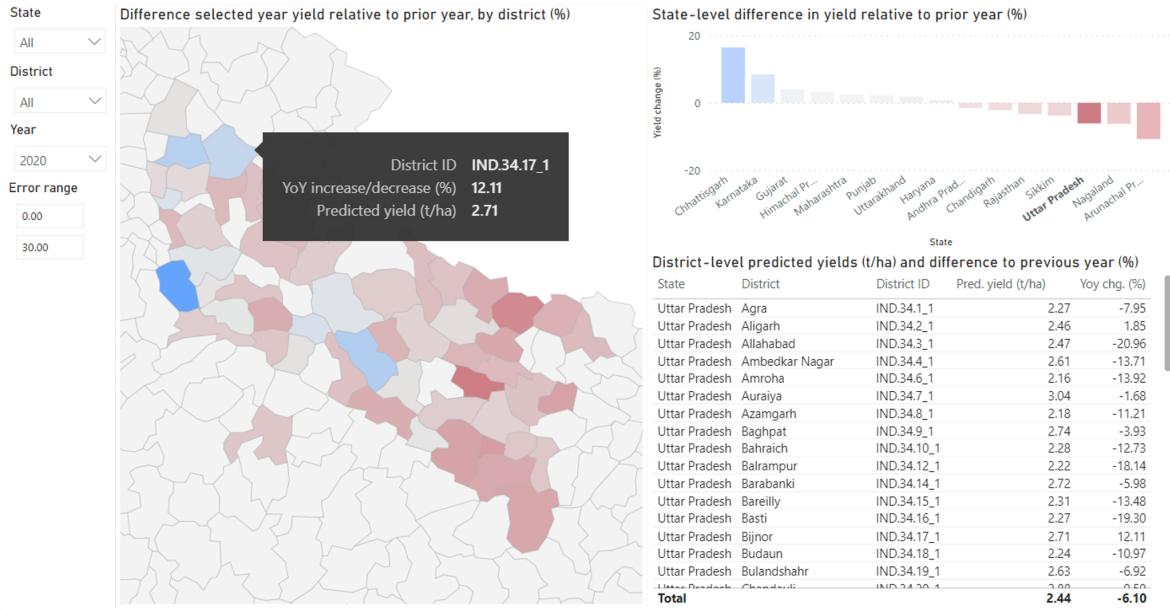
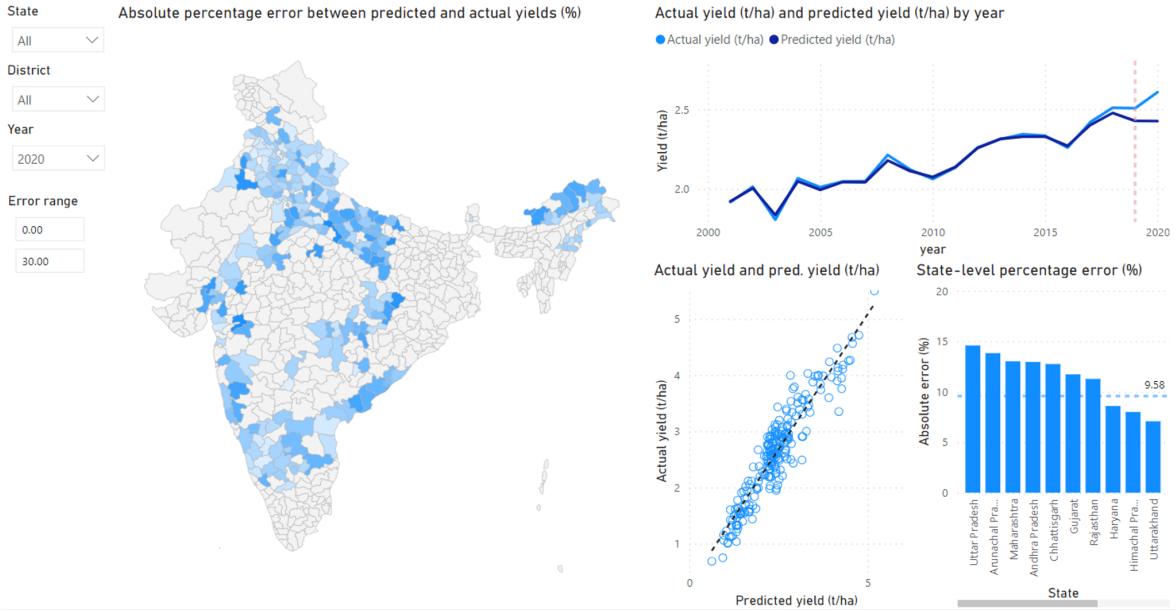


Figure 4: Interactive dashboard of yield prediction model outputs. The top panel shows a map of India with a colour coding applied to different districts to indicate the predicted yield values. Shades of blue indicate an increase in yield and shades of red denote a decrease in yield compared to the prior year's yield. Accompanying the map is a bar chart that provides a state-level summary and a table that enumerates the district-level predicted yields and the percentage change from the previous year across all states and districts. The bottom panel provides a similar comparative yield prediction, but focuses on the state of Uttar Pradesh.

INDIA RICE YIELD PREDICTION MODEL DIAGNOSTIC RESULTS



INDIA RICE YIELD PREDICTION MODEL DIAGNOSTIC RESULTS

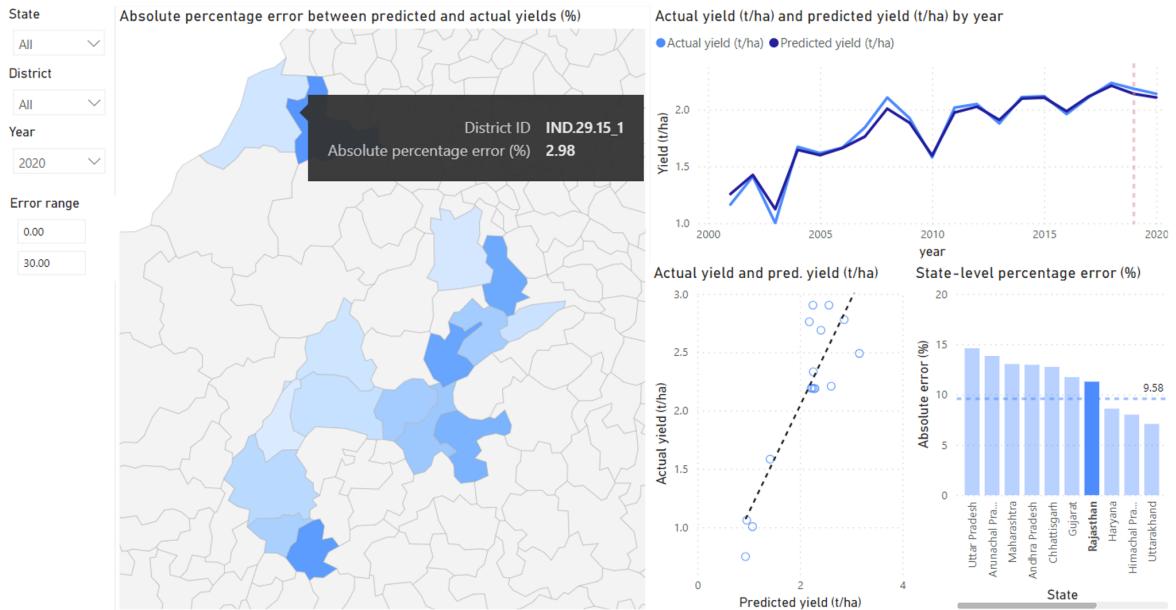


Figure 5: Interactive dashboard showing spatial view of yield prediction model error. The top panel provide model diagnostic information. The visual shows a map with the average percentage error by region; a scatter plot comparing the predicted yield and the actual yield; and a line graph showing the actual yield and predicted yield each year. The lower panel shows a similar view, zoomed in to districts within the state of Rajasthan.

including agronomy, suggests they could be effectively deployed in agricultural settings [88, 89]. For instance, existing digital solutions like KisanGPT, a chatbot designed to assist farmers with queries such as optimal fertilizer application, indicate the practical applications of LLMs in agriculture [90]. LLMs could help make rice yield prediction models more accessible to farmers through web platforms or mobile messaging applications, like WhatsApp, allowing for interaction in their native languages. Examples of mobile-based agricultural early warning systems have included wheat rust disease alerts in Ethiopia and weather alerts in Zimbabwe [91].

Yield prediction systems that leverage LLM technology to communicate outputs may also help bolster agricultural extension services. Bangladesh, for instance, has boosted support services for farmers, including enhanced information access and extension services, via the government's 'Info Sarkar' project, which aims to link government offices nationwide and has established over 4,500 internet-equipped Union Digital Centres (UDCs) to aid rural communities. Studies indicate these centers are being effectively utilized by educated youths, suggesting a potential for these individuals to lead in disseminating climate adaptation strategies to farmers [85].

Such agricultural extension services could benefit from additional digital aid such as LLMs that are able to distill agronomic science in a manner that is understandable and actionable for farmers, especially in contexts where agronomic advice should be disseminated in a context-specific manner. Advisories could be targeted based on whether farmers socio-economic status, the size of their farms, the level of internet connectivity, household income, amongst other factors. Past research has shown that differences in farmers backgrounds can significantly affect adoption of agricultural technologies, suggesting that tailored advice may be an important prerequisite for helping farmers adapt their agronomic practices in the face of climate-related risk [92].

Modelling the dissemination of yield prediction outputs. Building yield prediction systems that are tailored to local needs may help to boost resilience in agricultural communities. However, good design alone may be insufficient for widespread adoption. There is also a need to anticipate how external factors (such as mobile network coverage or regional agroclimatic variation) might affect the diffusion of this technology in a society across time and space, and how policymakers can enhance and sustain adoption.

The adoption of these systems could potentially be modeled using principles of the Diffusion of Innovation theory [93]. This theory provides a framework to understand how innovations are taken up in a population, highlighting the role of early adopters and the subsequent spread through social and communication networks. In the context of yield prediction systems, this might involve analyzing how the perceived advantages and compatibility with existing practices influence the rate of adoption among farmers or agricultural extension workers.

Furthermore, computational models of diffusion, such as the Bass diffusion model, can offer insights into the expected rate of technology uptake. These models provide a means to simulate and anticipate the adoption curve, taking into account various societal and technological factors [94]. By understanding the likely trajectory of technology diffusion, policies can be tailored to support and accelerate adoption, ensuring that the benefits of yield prediction systems are maximized. Applications of Bass-like models are numerous in the literature, but have yet to be applied to agricultural early warning technologies. Studies have applied the methodology to forecasting the diffusion of innovations including novel foods such as edible insects in the Netherland, groundwater pumps in Pakistan, preterm birth screening technology, fuel cell vehicles in China, household high-speed internet products, solar water heaters in Brazil, box office performance of movies, distributed solar generation, and cell phones in rural Bangladesh [95, 96, 97, 98, 99, 100, 101, 102, 103].

Leveraging yield prediction models for anticipatory aid allocation. Another potential area of research could involve investigating how to effectively integrate yield prediction algorithms into decision-making tools for the public sector to enhance the strategic allocation of humanitarian aid in response to weather-related challenges. For example, algorithms that produce anticipated rice yields in specific areas could be combined with optimization algorithms can assist government decisions on the anticipatory allocation of resources such as financial support, fertilizers, fungicides, and water resources among affected farming communities.

While examples of such optimization applied specifically to anticipatory aid allocation in agricultural settings are limited, there is a growing body of research on applications of operations research applied to humanitarian aid [104]. Examples of using optimization approaches (such as linear mixed integer programming, stochastic programming, or multi-objective optimization) in relevant contexts in-

clude: aid disbursement following the 2010 Haiti earthquake [105]; optimal aid disbursement in response to internal displacement in northwest Syria [106]; post-disaster distribution of essential humanitarian aid (medicine, food, and water) from temporary warehouses to points of demand in Peru [107]; optimal location planning of warehouse locations to store relief items in Thailand [108]; logistics distribution of essential relief items during COVID-19 lockdowns in Bangladesh [109]; optimization of United Nations Humanitarian Response Depot distribution plans [110]; humanitarian relief logistics in both pre- and post-disaster situations in presence of uncertainty [111]; and the World Food Programme’s Optimus tool, which leverages linear programming to optimize food aid operations [112, 113].

4 Conclusion

This study advances the state of the art in district-level rice yield prediction in India through an integrated approach, combining ERA5 climate reanalysis, MODIS satellite vegetation indices, and a novel, spatially matched yield dataset. By evaluating 19 machine learning models, the research establishes benchmarks for accuracy, achieving out-of-sample R², MAE, and MAPE values of up to 0.82, 0.29, and 0.16, respectively. The development of an interactive dashboard tool offers a means for visualizing yield predictions and assessing model performance across regions. This approach not only demonstrates the feasibility of using machine learning for rice yield forecasting at the district level in India, but also provides a benchmark for predictive accuracy against which further algorithmic innovations can be easily compared. The results offer evidence that machine learning-based rice yield prediction may have the potential to augment Indian agricultural early warning systems with robust crop yield prediction capabilities.

References

- [1] M. Weiss, F. Jacob, and G. Duveiller. Remote sensing for agricultural applications: a meta-review. *Remote Sens. Environ.*, 236:111402, January 2020.
- [2] Thomas van Klompenburg, Ayalew Kassahun, and Cagatay Catal. Crop yield prediction using machine learning: a systematic literature review. *Comput. Electron. Agric.*, 177:105709, October 2020.
- [3] Vasit Sagan, Maitiniyazi Maimaitijiang, Sourav Bhadra, Matthew Maimaitiyiming, Davis R. Brown, Paheding Sidike, and Felix B. Fritschi. Field-scale crop yield prediction using multi-temporal WorldView-3 and PlanetScope satellite data and deep learning. *ISPRS J. Photogramm. Remote Sens.*, 174:265–281, April 2021.
- [4] Ainong Li, Shunlin Liang, Angsheng Wang, and Jun Qin. Estimating crop yield from multi-temporal satellite data using multivariate regression and neural network techniques. *Photogramm. Eng. Remote Sens.*, 73(10):1149–1157, October 2007.
- [5] Ruixue Wang, Roderick M. Rejesus, and Serkan Aglasan. Warming temperatures, yield risk and crop insurance participation. *Eur. Rev. Agric. Econ.*, 48(5):1109–1131, December 2021.
- [6] Yvonne Wong Jing Wen, Raja Rajeswari Ponnusamy, and Ho Ming Kang. Application of weather index-based insurance for paddy yield: the case of Malaysia. *Int. J. Adv. Appl. Sci.*, 6:51–59, 2019.
- [7] Hari Sankar Nayak et al. Rice yield gaps and nitrogen-use efficiency in the Northwestern Indo-Gangetic Plains of India: evidence-based insights from heterogeneous farmers' practices. *Field Crops Research*, 275:108328, 2022.
- [8] FAO. India at a glance, 2018.
- [9] Muhammad Ishfaq et al. Alternate wetting and drying: a water-saving and ecofriendly rice production system. *Agricultural Water Management*, 241:106363, 2020.
- [10] Abdus Sattar and R. C. Srivastava. Modelling climate smart rice-wheat production system in the Middle Gangetic Plains of India. *Theoretical and Applied Climatology*, 144(1-2):77–91, 2021.
- [11] Steffen Fritz et al. A comparison of global agricultural monitoring systems and current gaps. *Agricultural Systems*, 168:258–272, 2019.
- [12] Xiangying Xu et al. Design of an integrated climatic assessment indicator (ICAI) for wheat production: a case study in Jiangsu Province, China. *Ecological Indicators*, 101:943–953, 2019.
- [13] Saeed Khaki, Lizhi Wang, and Sotirios V. Archontoulis. A CNN-RNN framework for crop yield prediction. *Frontiers in Plant Science*, 10:1750, 2020.
- [14] Anna X. Wang et al. Deep transfer learning for crop yield prediction with remote sensing data. In *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies*, 2018.
- [15] Hamid Kamangir, Brent S. Sams, Nick Dokoozlian, Luis Sanchez, and J. Mason Earles. Large-scale spatio-temporal yield estimation via deep learning using satellite and management data fusion in vineyards. *Computers and Electronics in Agriculture*, 216:108439, January 2024.
- [16] Phusanisa Charoen-Ung and Pradit Mittrapiyanuruk. Sugarcane yield grade prediction using random forest with forward feature selection and hyper-parameter tuning. In Herwig Unger, Sunantha Sodsee, and Phayung Meesad, editors, *Recent Advances in Information and Communication Technology 2018*, pages 33–42. Springer International Publishing, Cham, 2019.
- [17] Patrick Filippi et al. An approach to forecast grain crop yield using multi-layered, multi-farm data sets and machine learning. *Precision Agriculture*, 20(5):1015–1029, 2019.

- [18] Ishfaq Ahmad et al. Yield forecasting of spring maize using remote sensing and crop modeling in Faisalabad-Punjab Pakistan. *Journal of the Indian Society of Remote Sensing*, 46(10):–, 2018.
- [19] Anat Goldstein et al. Applying machine learning on sensor data for irrigation recommendations: revealing the agronomist's tacit knowledge. *Precision Agriculture*, 19(3):421–444, 2018.
- [20] Asheesh Chaurasiya et al. Layering smart management practices to sustainably maintain rice yields and improve water use efficiency in eastern India. *Field Crops Research*, 275:108341, 2022.
- [21] Zheng Chu and Jiong Yu. An end-to-end model for rice yield prediction using deep learning fusion. *Computers and Electronics in Agriculture*, 174:105471, 2020.
- [22] Li Tian et al. Yield prediction model of rice and wheat crops based on ecological distance algorithm. *Environmental Technology & Innovation*, 20:101132, 2020.
- [23] Weiguo Yu et al. Improved prediction of rice yield at field and county levels by synergistic use of SAR, optical and meteorological data. *Agricultural and Forest Meteorology*, 342:109729, 2023.
- [24] Liang Wan et al. Grain yield prediction of rice using multi-temporal UAV-based RGB and multispectral images and model transfer – a case study of small farmlands in the south of China. *Agricultural and Forest Meteorology*, 291:108096, 2020.
- [25] Prakash K Jha et al. Using daily data from seasonal forecasts in dynamic crop models for yield prediction: a case study for rice in Nepal's Terai. *Agricultural and Forest Meteorology*, 265:349–358, 2019.
- [26] Seungtaek Jeong, Jonghan Ko, and Jong-Min Yeom. Predicting rice yield at pixel scale through synthetic use of crop and deep learning models with satellite data in South and North Korea. *Science of The Total Environment*, 802:149726, 2022.
- [27] Ponraj Arumugam et al. Remote sensing based yield estimation of rice (*Oryza sativa L.*) using gradient boosted regression in India. *Remote Sensing*, 13(12):2379, 2021.
- [28] A Kumar Ranjan and B R Parida. Paddy acreage mapping and yield prediction using sentinel-based optical and SAR data in Sahibganj District, Jharkhand (India). *Spatial Information Research*, 27(4), 2019.
- [29] Diego Gómez et al. Regional estimation of garlic yield using crop, satellite and climate data in Mexico. *Computers and Electronics in Agriculture*, 181:105943, 2021.
- [30] Maximilian Auffhammer, V Ramanathan, and Jeffrey R Vincent. Climate change, the monsoon, and rice yield in India. *Climatic Change*, 111(2):411–424, 2012.
- [31] Sheetal Sharma et al. Field-specific nutrient management using rice crop manager decision support tool in Odisha, India. *Field Crops Research*, 241:107578, 2019.
- [32] Hans Hersbach et al. The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.
- [33] Aleš Urban et al. Evaluation of the ERA5 reanalysis-based universal thermal climate index on mortality data in Europe. *Environmental Research*, 198:111227, 2021.
- [34] Josh M Colston et al. Evaluating meteorological data from weather stations, and from satellites and global models for a multi-site epidemiological study. *Environmental Research*, 165:91–109, 2018.
- [35] Pramiti Kumar Chakraborty et al. Assessing congenial soil temperature and its impact on root growth, grain yield of summer rice under varying water stress condition in lower Gangetic Plain of India. *Journal of the Saudi Society of Agricultural Sciences*, 2021.
- [36] Y Jia et al. Effects of low water temperature during reproductive growth on photosynthetic production and nitrogen accumulation in rice. *Field Crops Research*, 242:107587, 2019.

- [37] Tsuneo Kuwagata et al. Hydrometeorology for plant omics: potential evaporation as a key index for transcriptome in rice. *Environmental and Experimental Botany*, page 104724, 2021.
- [38] Yongkang Tang et al. Effects of long-term low atmospheric pressure on gas exchange and growth of lettuce. *Advances in Space Research*, 46(6):751–760, 2010.
- [39] Ratneswar Poddar et al. Effect of irrigation regime and varietal selection on the yield, water productivity, energy indices and economics of rice production in the Lower Gangetic Plains of Eastern India. *Agricultural Water Management*, 2021.
- [40] Santanu Kumar Bal et al. Critical weather limits for paddy rice under diverse ecosystems of India. *Frontiers in Plant Science*, 14, 2023.
- [41] J. R. Alvarado. Influence of air temperature on rice population, length of period from sowing to flowering, and spikelet sterility. 2002.
- [42] Leah H Schinasi, Tarik Benmarhnia, and Anneclaire J De Roos. Modification of the association between high ambient temperature and health by urban microclimate indicators: a systematic review and meta-analysis. *Environmental Research*, 161:168–180, 2018.
- [43] Mariano F Lopresti, Carlos M Di Bella, and Américo J Degioanni. Relationship between MODIS-NDVI data and wheat yield: a case study in northern Buenos Aires Province, Argentina. *Information Processing in Agriculture*, 2(2):73–84, 2015.
- [44] Ewa Panek and Dariusz Gozdowski. Analysis of relationship between cereal yield and NDVI for selected regions of Central Europe based on MODIS satellite data. *Remote Sensing Applications: Society and Environment*, 17:100286, 2020.
- [45] N T Son et al. A comparative analysis of multitemporal MODIS EVI and NDVI data for large-scale rice yield estimation. *Agricultural and Forest Meteorology*, 197:52–64, 2014.
- [46] Fiona H M Tang et al. CROPGRIDS: a global geo-referenced dataset of 173 crops circa 2020. *Earth System Science Data Discussions*, pages 1–22, 2023.
- [47] Ministry of Agriculture and Farmers Welfare. Crop production statistics information system. Online, 2021.
- [48] Geetika Sonkar et al. Vulnerability of Indian wheat against rising temperature and aerosols. *Environmental Pollution*, 254:112946, 2019.
- [49] Chaitali Diwan et al. AI-based learning content generation and learning pathway augmentation to increase learner engagement. *Computers and Education: Artificial Intelligence*, 4:100110, 2023.
- [50] Global Administrative Areas. GADM database of global administrative areas, version 2.0. Online, 2012. Accessed on February 14, 2024.
- [51] Longfei Zhou et al. Improved yield prediction of ratoon rice using unmanned aerial vehicle-based multi-temporal feature method. *Rice Science*, 30(3):247–256, 2023.
- [52] Guolin Ke et al. LightGBM: a highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [53] Qi Shi, Mohamed Abdel-Aty, and Jaeyoung Lee. A Bayesian ridge regression analysis of congestion’s impact on urban expressway safety. *Accident Analysis & Prevention*, 88:124–137, 2016.
- [54] Michael E Tipping. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1:211–244, 2001.
- [55] David J C MacKay. Bayesian interpolation. In C Ray Smith, Gary J Erickson, and Paul O Neudorfer, editors, *Maximum Entropy and Bayesian Methods: Seattle, 1991*, pages 39–66. Springer Netherlands, 1992.

- [56] Helai Huang and Mohamed Abdel-Aty. Multilevel data and Bayesian analysis in traffic safety. *Accident Analysis & Prevention*, 42(6):1556–1565, 2010.
- [57] Jerome H Friedman. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378, 2002.
- [58] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [59] Peter Huber and Elvezio Ronchetti. *Robust statistics*. Wiley, 2009.
- [60] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning*. Springer New York, 2009.
- [61] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.
- [62] Trevor Hastie et al. Multi-class AdaBoost. *Statistics and Its Interface*, 2(3):349–360, 2009.
- [63] Ron Rubinstein, Michael Zibulevsky, and Michael Elad. Efficient implementation of the K-SVD algorithm using batch orthogonal matching pursuit. 2008.
- [64] Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*, 63(1):3–42, 2006.
- [65] Vitor Cerqueira, Luis Torgo, and Igor Mozetič. Evaluating time series forecasting models: An empirical study on performance estimation methods. *Machine Learning*, 109(11):1997–2028, November 2020.
- [66] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, pages 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [67] Ian Lenaers and Lieven De Moor. Exploring XAI techniques for enhancing model transparency and interpretability in real estate rent prediction: A comparative study. *Finance Research Letters*, 58:104306, December 2023.
- [68] Vipul Jain, H Kavitha, and S Mohana Kumar. Credit card fraud detection web application using Streamlit and machine learning. In *2022 IEEE International Conference on Data Science and Information System (ICDSIS)*, pages 1–5, 2022.
- [69] Chanin Nantasenamat and et al. Chapter 27 - Building bioinformatics web applications with Streamlit. In Kunal Roy, editor, *Cheminformatics, QSAR and Machine Learning Applications for Novel Drug Development*, pages 679–699. Academic Press, 2023.
- [70] Shilpa Patil and V. Lokesha. Live Twitter sentiment analysis using Streamlit framework. *SSRN Scholarly Paper*, May 2022.
- [71] Mariem Belghith and et al. A new rolling forecasting framework using Microsoft Power BI for data visualization: a case study in a pharmaceutical industry. *Annales Pharmaceutiques Françaises*, November 2023.
- [72] Kathy Abusager, Michael Baldwin, and Vincent Hsu. Using Power BI to inform Clostridiooides difficile ordering practices at an acute care hospital in Central Florida. *American Journal of Infection Control*, 48(8, Supplement):S57–S58, August 2020.
- [73] Erin Burrell Nickell, Jason Schwebke, and Paul Goldwater. An introductory audit data analytics case study: using Microsoft Power BI and Benford’s law to detect accounting irregularities. *Journal of Accounting Education*, 64:100855, September 2023.
- [74] Yahui Guo et al. Integrated phenology and climate in rice yields prediction using machine learning methods. *Ecological Indicators*, 120:106935, 2021.

- [75] Kersten Clauss et al. Estimating rice production in the Mekong Delta, Vietnam, utilizing time series of Sentinel-1 SAR data. *International Journal of Applied Earth Observation and Geoinformation*, 73:574–585, 2018.
- [76] Jingye Han et al. Rice yield estimation using a CNN-based image-driven data assimilation framework. *Field Crops Research*, 288:108693, 2022.
- [77] Xi Su et al. Grain yield prediction using multi-temporal UAV-based multispectral vegetation indices and endmember abundance in rice. *Field Crops Research*, 299:108992, 2023.
- [78] Md. Monirul Islam et al. Development of remote sensing-based yield prediction models at the maturity stage of Boro rice using parametric and nonparametric approaches. *Remote Sensing Applications: Society and Environment*, 22:100494, 2021.
- [79] X Zhou et al. Predicting grain yield in rice using multi-temporal vegetation indices from UAV-based multispectral and digital imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 130:246–255, 2017.
- [80] Tri D Setiyono et al. Spatial rice yield estimation based on MODIS and Sentinel-1 SAR data and ORYZA crop growth model. *Remote Sensing*, 10(2), 2018.
- [81] Senlin Guan et al. Assessing correlation of high-resolution NDVI with fertilizer application level and yield of rice and wheat crops using small UAVs. *Remote Sensing*, 11(2), 2019.
- [82] Vitor Cerqueira, Nuno Moniz, and Carlos Soares. VEST: automatic feature engineering for forecasting. *Machine Learning*, 2021.
- [83] Jiaping Liang et al. Analysis and prediction of the impact of socio-economic and meteorological factors on rapeseed yield based on machine learning. *Agronomy*, 13(7), 2023.
- [84] Roel Jongeneel and Ana Rosa Gonzalez-Martinez. Estimating crop yield supply responses to be used for market outlook models: application to major developed and developing countries. *NJAS - Wageningen Journal of Life Sciences*, 92, 2020.
- [85] Zeenatul Islam, Mohammad Alauddin, and Md. Abdur Rashid Sarker. Determinants and implications of crop production loss: an empirical exploration using ordered probit analysis. *Land Use Policy*, 67:527–536, 2017.
- [86] Vipul Mann et al. SUSIE: pharmaceutical CMC ontology-based information extraction for drug development using machine learning. *Computers & Chemical Engineering*, 179:108446, 2023.
- [87] Han Zhang et al. Towards foundation models for learning on tabular data. *arXiv*, 2023.
- [88] Bernardino Romera-Paredes et al. Mathematical discoveries from program search with large language models. *Nature*, pages 1–3, 2023.
- [89] Bruno Silva et al. GPT-4 as an agronomist assistant? Answering agriculture exams using large language models. *arXiv*, 2023.
- [90] A. Tzachor et al. Large language models and agricultural extension services. *Nature Food*, 4(11):941–948, 2023.
- [91] Clare Allen-Sader et al. An early warning system to predict and mitigate wheat rust diseases in Ethiopia. *Environmental Research Letters*, 14(11):115004, 2019.
- [92] Parshuram Samal et al. Rice ecosystems and factors affecting varietal adoption in rainfed coastal Orissa: a multivariate probit analysis. *Agricultural Economics Research Review*, 24(1), 2011.
- [93] Everett M. Rogers. *Diffusion of innovations*. Simon and Schuster, 2010.
- [94] Frank M. Bass. A new product growth for model consumer durables. *Management Science*, 15(5):215–227, 1969.

- [95] Andrijana Horvat, Vincenzo Fogliano, and Pieterneel A. Luning. Modifying the Bass diffusion model to study adoption of radical new foods—The case of edible insects in the Netherlands. *PLOS ONE*, 15(6):e0234538, 2020.
- [96] A. Siddiqi, J. Tran, and J. L. Jr. Wescoat. Modeling growth and diffusion of groundwater pumping at multiple sub-provincial scales. In *AGU Fall Meeting Abstracts*, volume 2018, pages GC51I–0906, 2018.
- [97] Sabine E. Grimm, John W. Stevens, and Simon Dixon. Estimating future health technology diffusion using expert beliefs calibrated to an established diffusion model. *Value in Health*, 21(8):944–950, 2018.
- [98] Y. Xian et al. Research on the market diffusion of fuel cell vehicles in China based on the generalized Bass model. *IEEE Transactions on Industry Applications*, 58(2):2950–2960, 2022.
- [99] Franklin M. Lartey. Predicting product uptake using Bass, Gompertz, and Logistic diffusion models: application to a broadband product. *Journal of Business Administration Research*, 9(2):5, 2020.
- [100] Hendrigo Batista da Silva, Wadaed Uturbey, and Bruno M. Lopes. Market diffusion of household PV systems: Insights using the Bass model and solar water heaters market data. *Energy for Sustainable Development*, 55:210–220, 2020.
- [101] Chuan Zhang, Yu-Xin Tian, and Zhi-Ping Fan. Forecasting the box offices of movies coming soon using social media analysis: A method based on improved Bass models. *Expert Systems with Applications*, 191:116241, 2022.
- [102] Ryan Ratcliff and Kokila Doshi. Using the Bass model to analyze the diffusion of innovations at the base of the pyramid. *Business & Society*, 55(2):271–298, 2016.
- [103] Tiago P. Abud et al. A modified Bass model to calculate PVDG hosting capacity in LV networks. *Electric Power Systems Research*, 209:107966, 2022.
- [104] Walter J. Gutjahr and Pamela C. Nolz. Multicriteria optimization in humanitarian aid. *European Journal of Operational Research*, 252(2):351–366, 2016.
- [105] Begoña Vitoriano et al. A multi-criteria optimization model for humanitarian aid distribution. *Journal of Global Optimization*, 51(2):189–208, 2011.
- [106] Israa Ismail. A probabilistic mathematical programming model to control the flow of relief commodities in humanitarian supply chains. *Computers & Industrial Engineering*, 157:107305, 2021.
- [107] Alvaro Alonso Acero Condor, Cesar Manuel Ramirez Castañeda, and José Antonio Taquía Gutiérrez. Optimization of humanitarian aid resource distribution time through mixed integer linear programming. In *Proceedings of the 2023 9th International Conference on Industrial and Business Engineering*, pages 191–197, New York, NY, USA, 2023. Association for Computing Machinery.
- [108] Chawis Boonmee and Chompoonoot Kasemset. The multi-objective fuzzy mathematical programming model for humanitarian relief logistics. *Industrial Engineering & Management Systems*, 19(1):197–210, 2020.
- [109] Ziaul Haq Adnan et al. Applying linear programming for logistics distribution of essential relief items during COVID-19 lockdown: Evidence from Bangladesh. *International Journal of Logistics Economics and Globalisation*, 9(3):191–204, 2022.
- [110] İbrahim Miraç Eligüzel, Eren Özceylan, and Gerhard-Wilhelm Weber. Location-allocation analysis of humanitarian distribution plans: A case of United Nations Humanitarian Response Depots. *Annals of Operations Research*, 324(1):825–854, 2023.
- [111] Peiman Ghasemi, Fariba Goodarzian, and Ajith Abraham. A new humanitarian relief logistic network for multi-objective optimization under stochastic programming. *Applied Intelligence*, 52(12):13729–13762, 2022.

- [112] Koen Peters et al. The nutritious supply chain: Optimizing humanitarian food assistance. *INFORMS Journal on Optimization*, 3(2):200–226, 2021.
- [113] Koen Peters et al. UN World Food Programme: Toward zero hunger with analytics. *INFORMS Journal on Applied Analytics*, 52(1):8–26, 2022.