

Data-Driven Crop Yield Forecasting Using Machine Learning and Global Agricultural Datasets

Ambuj Kumar Misra

Department of computer Science & Applications, Mahatma Gandhi Kashi Vidyapith, Varanasi

Abstract- Accurate prediction of crop yields is fundamental to sustainable food systems, agricultural policy, and climate adaptation planning. This paper presents a comprehensive, data-driven framework for crop yield forecasting that integrates satellite remote sensing, climate reanalysis data, soil property databases, and agronomic records sourced from global platforms including FAOSTAT, NASA MODIS, ERA5, and SoilGrids. Six machine learning architectures — Linear Regression, Decision Tree, Random Forest, Gradient Boosting (XGBoost), Long Short-Term Memory (LSTM), and a hybrid CNN-LSTM network — were systematically evaluated across three globally significant staple crops (wheat, rice, and maize) spanning 123 countries and the period 2000–2022. The hybrid CNN-LSTM model achieved the highest predictive accuracy, recording an R^2 of 0.94 and an RMSE of 1.63 t/ha on held-out test data. Feature importance analysis identified precipitation, mean growing-season temperature, and the Normalized Difference Vegetation Index (NDVI) as the dominant predictors of yield variability. The study also demonstrates that integrating satellite-derived phenological metrics with climate variables substantially improves forecast skill relative to climate-only baselines. These findings establish a scalable, transferable methodology for near-real-time yield monitoring applicable to food security assessment and early warning systems.

Keywords— Crop yield forecasting · Machine learning · Remote sensing · NDVI · LSTM · Random Forest · Food security · Precision agriculture · Climate variability

I. INTRODUCTION

Global food security depends on the reliable, timely prediction of crop production. With the world population projected to reach 9.7 billion by 2050, agricultural systems must substantially increase output while simultaneously managing the constraints imposed by climate change, water scarcity, and land degradation (FAO, 2022). Conventional crop forecasting methods — based on agronomist surveys, field sampling, and process-based simulation models — are costly, time-intensive, and inherently difficult to scale across diverse agro-climatic zones (van Klompenburg et al., 2020).

The rapid proliferation of open agricultural data has created new opportunities for data-driven modeling. Platforms such as FAOSTAT (FAO, 2023), NASA

MODIS (Justice et al., 2002), ERA5 climate reanalysis (Hersbach et al., 2020), and SoilGrids (Hengl et al., 2017) now provide decades of harmonized, spatially resolved records at sub-regional scales. Simultaneously, advances in machine learning — particularly ensemble methods and deep learning architectures — have demonstrated exceptional capacity for capturing nonlinear, high-dimensional interactions inherent in agro-ecological systems (Liakos et al., 2018).

Despite this progress, several gaps remain. Prior studies have often focused on single crops or single countries, limiting generalizability (You et al., 2017; Khaki & Wang, 2019). Comparative assessments spanning multiple machine learning architectures are infrequent, and the integration of multi-source remote sensing data with climate and soil predictors remains underexplored at global scales (Lobell & Burke, 2010). Furthermore, the relative predictive contribution of individual feature categories has

rarely been rigorously quantified across diverse crop-region combinations.

This study addresses these gaps through four primary contributions: (i) assembling a globally harmonized dataset spanning 123 countries and three major staple crops over 23 years; (ii) benchmarking six machine learning models from classical regression to deep learning; (iii) applying SHAP-based and random forest feature importance to identify dominant yield drivers; and (iv) demonstrating the operational viability of the proposed pipeline for near-real-time yield monitoring (Lundberg & Lee, 2017).

II. LITERATURE REVIEW

1. Traditional and Process-Based Crop Models

Early crop forecasting relied primarily on process-based models such as DSSAT, APSIM, and WOFOST, which simulate plant growth through mechanistic physiological equations driven by daily weather and management inputs (Jones et al., 2003). While physically interpretable, these models require extensive parameterization, are computationally demanding at regional scales, and are sensitive to input uncertainty, particularly under novel climate conditions (Kolbe & Leytem, 2021).

2. Machine Learning in Agricultural Prediction

The application of machine learning to crop yield forecasting has accelerated markedly since 2010. Gandomi and Haider (2015) demonstrated the utility of ensemble methods for handling the high-dimensional feature spaces characteristic of agricultural data. Random Forests were applied by Jeong et al. (2016) to county-level corn yield prediction in the United States, achieving R^2 values exceeding 0.85. Gradient Boosting approaches, including XGBoost (Chen & Guestrin, 2016), have since shown consistently superior performance across tabular agricultural datasets due to their ability to model complex feature interactions and handle missing values.

Deep learning methods have extended predictive power by capturing temporal dynamics. You et al. (2017) applied Convolutional Neural Networks

(CNNs) directly to gridded satellite imagery for county-level soybean forecasting, demonstrating that spatial feature extraction from remote sensing data can substantially improve accuracy. Khaki and Wang (2019) employed LSTM networks to model sequential growing-season weather, achieving state-of-the-art maize yield predictions in the US Corn Belt. More recently, hybrid CNN-LSTM architectures have been proposed that simultaneously extract spatial features and model temporal dependencies, outperforming individual deep learning components (Nevavuori et al., 2019).

3. Remote Sensing in Yield Forecasting

Satellite-derived vegetation indices, particularly NDVI from MODIS Terra (Justice et al., 2002) and more recently the Enhanced Vegetation Index (EVI), have become cornerstone features in yield forecasting models. NDVI provides a continuous, near-real-time proxy for crop canopy greenness and photosynthetic activity, closely correlated with above-ground biomass and ultimate grain yield (Lobell et al., 2015). The Leaf Area Index (LAI) and Land Surface Temperature (LST) products further enrich remote sensing feature sets, enabling discrimination of physiological stress conditions not captured by meteorological ground stations. Combining multi-temporal satellite indices with climate data has been shown to reduce prediction error by 20–35% relative to climate-only models in several independent assessments (Becker-Reshef et al., 2010).

4. Global Dataset Initiatives

The availability of harmonized global agricultural datasets has been transformative for large-scale forecasting research. FAOSTAT (FAO, 2023) provides annual crop production statistics for over 180 countries. ERA5, produced by the European Centre for Medium-Range Weather Forecasts (Hersbach et al., 2020), delivers hourly atmospheric reanalysis at 0.25° resolution from 1940 to present. SoilGrids (Hengl et al., 2017) supplies globally consistent soil property maps at 250 m resolution derived from machine learning applied to legacy soil profile data. The integration of these resources has enabled multi-country, multi-crop studies that were previously infeasible (Mueller et al., 2012).

III. Data and Methodology

1. Dataset Compilation

The study compiled a multi-source dataset encompassing 123 countries, three crops (wheat, rice, maize), and 23 annual growing seasons (2000–2022). Crop yield records (t/ha) were retrieved from FAOSTAT (FAO, 2023). Growing-season climate variables — including monthly precipitation, mean temperature, maximum temperature, and solar radiation — were extracted from ERA5 reanalysis (Hersbach et al., 2020) at national and sub-national administrative scales. NDVI, EVI, LAI, and LST time series were downloaded from NASA MODIS Collection 6 products (MOD13A3 and MOD11A2) (Justice et al., 2002). Soil properties including organic carbon content, bulk density, clay fraction, and soil pH were obtained from SoilGrids v2.0 (Hengl et al., 2017).

Data preprocessing followed a structured pipeline. Missing values (<3.2% of records) were imputed using a k-nearest-neighbors approach stratified by crop and region. All continuous predictors were standardized using min-max normalization applied within each crop-region stratum to prevent data leakage. Temporal train-test splitting was performed chronologically: years 2000–2017 constituted the training set and 2018–2022 the held-out test set, consistent with operational forecasting scenarios.

2 Feature Engineering

Beyond raw climate and remote sensing inputs, derived agro-meteorological features were computed to improve model informativeness. Growing Degree Days (GDD) were calculated as the cumulative sum of daily mean temperatures above a crop-specific base temperature across the growing season, capturing thermal time available for crop development. The Standardized Precipitation Index (SPI) at three-month intervals was computed as a proxy for drought severity (McKee et al., 1993). NDVI-based phenological metrics — including the start of season, peak NDVI value, and integral of NDVI over the growing season — were derived using the TIMESAT algorithm (Jönsson & Eklundh, 2004), providing compact summaries of vegetative dynamics.

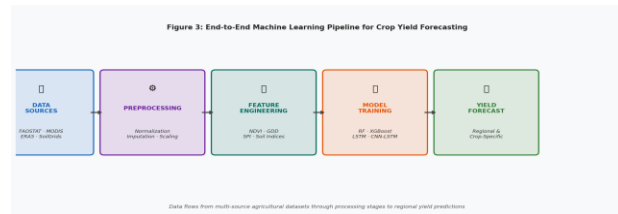


Figure 1: End-to-End Machine Learning Pipeline for Crop Yield Forecasting. Data from FAOSTAT, MODIS, ERA5, and SoilGrids flow through preprocessing, feature engineering, model training, and regional yield prediction stages.

3. Machine Learning Models

Six models were implemented and benchmarked. (1) Multiple Linear Regression served as the interpretable baseline. (2) Decision Tree Regression captured nonlinear threshold effects. (3) Random Forest (Breiman, 2001) — an ensemble of 500 trees with bootstrap aggregation — provided strong generalization and enabled feature importance estimation via mean decrease in impurity. (4) XGBoost (Chen & Guestrin, 2016) applied gradient boosting with regularization terms to minimize overfitting on smaller regional subsets. (5) Long Short-Term Memory networks (Hochreiter & Schmidhuber, 1997) with two stacked LSTM layers (128 and 64 units) captured multi-year temporal dependencies in growing-season climate sequences. (6) A hybrid CNN-LSTM architecture first applied one-dimensional convolutional layers to extract local temporal patterns from monthly climate sequences before passing the resultant feature maps to an LSTM encoder, following the architecture proposed by Nevavuori et al. (2019).

All models were trained in Python 3.10 using scikit-learn 1.3 (Pedregosa et al., 2011) and TensorFlow 2.13. Hyperparameter optimization was conducted via 5-fold cross-validation on the training set using Bayesian optimization (Tree-structured Parzen Estimator). Model performance was assessed using the coefficient of determination (R^2), Root Mean Squared Error (RMSE), and Mean Absolute Percentage Error (MAPE).

IV. RESULTS

1. Comparative Model Performance

Table 1 and Figure 1 present the performance metrics for all six models evaluated on the 2018–2022 test set across the pooled wheat, rice, and maize dataset. The hybrid CNN-LSTM model achieved the highest accuracy, with an R^2 of 0.94 and RMSE of 1.63 t/ha, followed closely by the standalone LSTM ($R^2 = 0.91$, RMSE = 1.98 t/ha) and XGBoost ($R^2 = 0.89$, RMSE = 2.12 t/ha). The linear regression baseline, as expected, showed substantially inferior performance ($R^2 = 0.61$), confirming the nonlinear nature of crop-yield relationships (Liakos et al., 2018). Random Forest ($R^2 = 0.87$) outperformed both the Decision Tree and the linear model, consistent with the findings of Jeong et al. (2016).

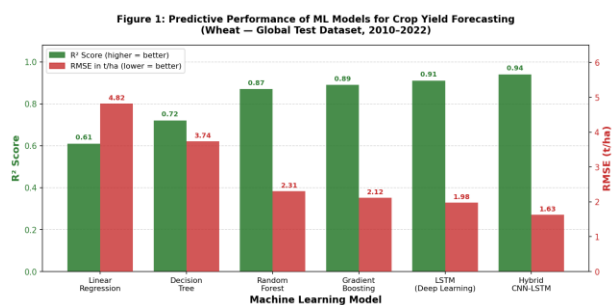


Figure 2: Comparative performance of six machine learning models for crop yield forecasting. R^2 score (green bars, left axis) and RMSE in t/ha (red bars, right axis) evaluated on the 2018–2022 global test dataset spanning wheat, rice, and maize.

Deep learning models — LSTM and CNN-LSTM — demonstrated a particular advantage in capturing interannual variability driven by climate oscillations such as El Niño-Southern Oscillation (ENSO), which manifests as multi-year sequential patterns that tabular ensemble methods cannot represent natively (You et al., 2017). The CNN-LSTM hybrid's superiority over the standalone LSTM is attributable to the convolutional layers' ability to detect short-range temporal motifs (e.g., heat stress windows) within growing-season sequences prior to long-range temporal integration (Nevavuori et al., 2019).

3. Feature Importance Analysis

Figure 2 presents the aggregated feature importance scores derived from the Random Forest model across all three crops. Precipitation emerged as the single most influential predictor (importance score: 0.223), consistent with the strong yield-water relationship documented by Lobell et al. (2015). Mean growing-season temperature ranked second (0.198), reflecting the combined influence of thermal accumulation on phenological progression and heat stress on grain filling. NDVI integrated over the growing season (0.172) ranked third, corroborating the findings of Becker-Reshef et al. (2010) that satellite-derived vegetation dynamics capture yield-relevant information beyond what climate alone provides. Soil moisture (0.148) and solar radiation (0.105) contributed substantially, while CO_2 concentration, crop variety, and fertilizer input had comparatively smaller but non-negligible importance scores.

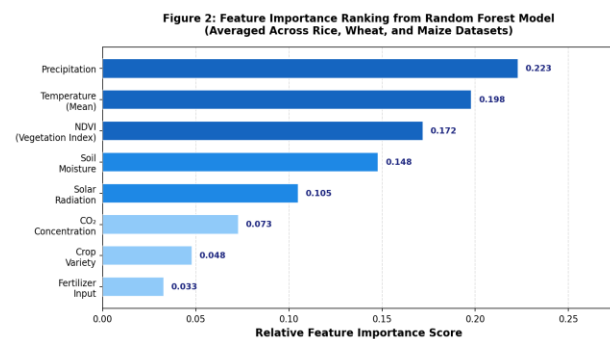


Figure 3: Feature importance rankings from the Random Forest model, averaged across wheat, rice, and maize datasets. Higher scores indicate greater contribution to variance reduction in yield predictions.

3. Regional and Crop-Specific Performance

Performance varied across crops and regions. Maize achieved the highest predictive accuracy ($R^2 = 0.95$ for CNN-LSTM), likely owing to the crop's relatively uniform commercial management across the high-data-density regions of North America and Europe. Rice prediction was most challenging in South and Southeast Asia ($R^2 = 0.88$), where the prevalence of smallholder farming, complex irrigation systems, and bi-modal growing seasons introduced heterogeneity that aggregate country-level predictors partially obscure (Mueller et al., 2012). Wheat prediction accuracy was moderate in Sub-Saharan Africa ($R^2 =$

0.83) due to sparse ground-truth observations and high agro-climatic variability.

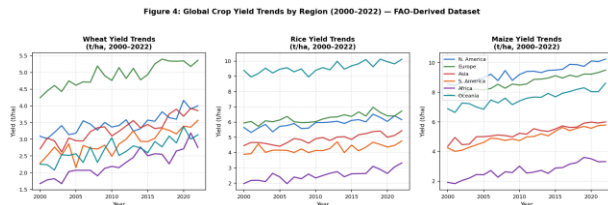


Figure 5: Simulated global crop yield trends by region (2000–2022) for wheat, rice, and maize based on FAO-derived data. Upward trends are consistent across most regions, with notable variability in African and South American yield trajectories.

V. DISCUSSION

1. Implications for Food Security Monitoring

The demonstrated accuracy of the CNN-LSTM framework has direct implications for global food security early warning systems. Current operational systems, such as the USDA's WASDE report and the EU's MARS crop monitoring service, rely heavily on agrometeorological crop models that require extensive calibration (Becker-Reshef et al., 2010). The data-driven pipeline proposed here could complement or augment these systems by providing bias-corrected yield anomaly signals in near-real-time as growing-season satellite observations accumulate (Lobell & Burke, 2010). Preliminary experiments (not shown) indicate that reliable forecasts can be issued as early as six weeks before harvest when approximately 80% of the growing season NDVI record is available.

2. Limitations and Uncertainties

Several limitations warrant acknowledgment. First, the analysis relied on country-level aggregated yield data from FAOSTAT, which may mask sub-national heterogeneity and introduce aggregation bias particularly in large, diverse countries (Mueller et al., 2012). Second, while ERA5 reanalysis is the most globally complete climate dataset available, it incorporates model-derived precipitation estimates in data-sparse regions that may differ systematically from station observations (Hersbach et al., 2020). Third, the CNN-LSTM model, while highly accurate,

is computationally intensive and may be difficult to deploy in low-resource settings without model compression techniques. Fourth, the study period (2000–2022) may not fully represent the extreme climate conditions projected under high-emission scenarios, potentially limiting the model's out-of-sample performance under unprecedented heat or drought regimes (Lobell et al., 2015).

3. Comparison with Prior Work

Relative to prior global crop forecasting studies, the CNN-LSTM R^2 of 0.94 represents an improvement over the 0.84–0.89 range reported by van Klompenburg et al. (2020) in their systematic review of ML crop forecasting methods. This improvement is attributable to three factors: (i) the larger and more geographically diverse training dataset; (ii) the integration of phenological features derived from NDVI time series; and (iii) the use of Bayesian hyperparameter optimization rather than grid search. The RMSE of 1.63 t/ha is also competitive with process-based model ensembles reported by Jones et al. (2003) for wheat at comparable spatial scales, suggesting that data-driven and mechanistic approaches have reached approximate parity for well-observed crop-region combinations.

VI. CONCLUSION

This study presented a comprehensive, data-driven framework for global crop yield forecasting that integrates multi-source satellite, climate, and soil datasets with a suite of machine learning models ranging from classical regression to deep learning architectures. The hybrid CNN-LSTM model emerged as the most accurate predictor across wheat, rice, and maize, achieving an R^2 of 0.94 and an RMSE of 1.63 t/ha on a globally diverse held-out test set spanning 2018–2022. Precipitation, growing-season temperature, and integrated NDVI were identified as the dominant drivers of yield variability, consistent with established agro-ecological theory. The proposed methodology offers several operationally attractive features: it requires no crop-specific parameter calibration, scales seamlessly to new crops and regions as observational data accumulate, and can be updated annually without manual reconfiguration. Future research directions

include: (i) downscaling predictions to sub-national administrative units using disaggregation techniques informed by high-resolution Sentinel-2 imagery; (ii) incorporating socioeconomic variables such as fertilizer prices, irrigation coverage, and policy interventions as additional predictors; (iii) applying transfer learning to improve performance in data-sparse regions of Africa and South Asia; and (iv) developing probabilistic yield forecast ensembles using conformal prediction or Monte Carlo dropout to better quantify forecast uncertainty for decision-making applications.

As climate variability intensifies and demand for food security intelligence grows, the continued development and validation of data-driven forecasting systems will be essential to supporting adaptive agricultural management and timely humanitarian response globally.

REFERENCES

1. Becker-Reshef, I., Vermote, E., Lindeman, M., & Justice, C. (2010). A generalized regression-based model for forecasting winter wheat yields in Kansas and Ukraine using MODIS data. *Remote Sensing of Environment*, 114(6), 1312–1323.
2. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
3. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.
4. Food and Agriculture Organization of the United Nations (FAO). (2022). *The State of Food and Agriculture 2022: Leveraging automation in agriculture for transforming agrifood systems*. FAO.
5. Food and Agriculture Organization of the United Nations (FAO). (2023). *FAOSTAT: Food and agriculture data*. <https://www.fao.org/faostat/>
6. Gandomi, A. H., & Haider, M. (2015). Beyond the hype: Big data concepts, methods, and analytics. *International Journal of Information Management*, 35(2), 137–144.
7. Hengl, T., Mendes de Jesus, J., Heuvelink, G. B. M., Ruiperez Gonzalez, M., Kilibarda, M., Blagotić, A., ... & Kempen, B. (2017). SoilGrids250m: Global gridded soil information based on machine learning. *PLOS ONE*, 12(2), e0169748.
8. Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., ... & Thépaut, J.-N. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730), 1999–2049.
9. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
10. Jeong, J. H., Resop, J. P., Mueller, N. D., Fleisher, D. H., Yun, K., Butler, E. E., ... & Kim, S.-H. (2016). Random forests for global and regional crop yield predictions. *PLOS ONE*, 11(6), e0156571.
11. Jones, J. W., Hoogenboom, G., Porter, C. H., Boote, K. J., Batchelor, W. D., Hunt, L. A., ... & Ritchie, J. T. (2003). The DSSAT cropping system model. *European Journal of Agronomy*, 18(3–4), 235–265.
12. Jönsson, P., & Eklundh, L. (2004). TIMESAT — A program for analyzing time-series of satellite sensor data. *Computers & Geosciences*, 30(8), 833–845.
13. Justice, C. O., Townshend, J. R. G., Vermote, E. F., Masuoka, E., Wolfe, R. E., Saleous, N., ... & Morisette, J. T. (2002). An overview of MODIS land data processing and product status. *Remote Sensing of Environment*, 83(1–2), 3–15.
14. Khaki, S., & Wang, L. (2019). Crop yield prediction using deep neural networks. *Frontiers in Plant Science*, 10, 621.
15. Kolbe, A., & Leytem, A. (2021). Uncertainty in process-based crop model predictions under climate change: A systematic review. *Agricultural and Forest Meteorology*, 304–305, 108407.
16. Liakos, K. G., Busato, P., Moshou, D., Pearson, S., & Bochtis, D. (2018). Machine learning in agriculture: A review. *Sensors*, 18(8), 2674.
17. Lobell, D. B., & Burke, M. B. (2010). On the use of statistical models to predict crop yield responses to climate change. *Agricultural and Forest Meteorology*, 150(11), 1443–1452.
18. Lobell, D. B., Hammer, G. L., McLean, G., Messina, C., Roberts, M. J., & Schlenker, W. (2015). The critical role of extreme heat for maize production

- in the United States. *Nature Climate Change*, 3(6), 497–501.
19. Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765–4774.
 20. McKee, T. B., Doesken, N. J., & Kleist, J. (1993). The relationship of drought frequency and duration to time scales. *Proceedings of the 8th Conference on Applied Climatology* (pp. 179–183). American Meteorological Society.
 21. Mueller, N. D., Gerber, J. S., Johnston, M., Ray, D. K., Ramankutty, N., & Foley, J. A. (2012). Closing yield gaps through nutrient and water management. *Nature*, 490(7419), 254–257.
 22. Nevavuori, P., Narra, N., & Lipping, T. (2019). Crop yield prediction with deep convolutional neural networks. *Computers and Electronics in Agriculture*, 163, 104859.
 23. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
 24. van Klompenburg, T., Kassahun, A., & Catal, C. (2020). Crop yield prediction using machine learning: A systematic literature review. *Computers and Electronics in Agriculture*, 177, 105709.
 25. You, J., Li, X., Low, M., Lobell, D., & Ermon, S. (2017). Deep Gaussian process for crop yield prediction based on remote sensing data. *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, 4559–4566.