



Data-Driven Agricultural Decision Support System

Ambuj Kumar Misra

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

Abstract: Contemporary agriculture operates under the compounding pressures of climate variability, resource scarcity, and an expanding global population. This paper introduces a comprehensive Data-Driven Agricultural Decision Support System (DADSS) that harnesses decades of historical crop production records, real-time sensor telemetry, and satellite-derived remote sensing imagery to generate actionable, site-specific recommendations for farm management. Four machine learning architectures — Random Forest (RF), Long Short-Term Memory (LSTM) networks, Gradient Boosting Machines (GBM), and Support Vector Machines (SVM) — were trained and benchmarked against a baseline linear regression model across six major U.S. cropping regions spanning the years 2000 to 2022. The hybrid DADSS framework, which fuses LSTM temporal modeling with gradient-boosted ensemble predictions, achieved an overall forecasting accuracy of 92.5%, outperforming all individual baselines. Results confirm that data-driven advisory tools can meaningfully reduce input costs, improve yield stability, and bolster farmers' adaptive capacity in the face of climate uncertainty. The system architecture, feature engineering pipeline, validation results, and policy implications are discussed in detail.

Keywords: Agricultural Decision Support System · Machine Learning · Crop Yield Prediction · Precision Agriculture · Data Analytics · Historical Crop Data · Random Forest · LSTM Neural Network

I. Introduction

Agriculture remains the bedrock of global food security, providing livelihoods for more than 1.3 billion people worldwide and accounting for roughly 4% of global GDP. Yet the sector confronts an unprecedented convergence of threats: erratic precipitation and rising temperatures linked to climate change, mounting competition for freshwater resources, declining soil organic matter, and volatile commodity markets that punish uninformed planting decisions. [3] [19]

Traditional farm management has long relied on accumulated local knowledge, intuition, and rudimentary trial-and-error experimentation. While this inherited wisdom carries genuine agronomic value, it is inherently limited in its ability to synthesize the multidimensional, nonlinear relationships among soil properties, weather dynamics, pest pressures, and genetic performance of modern cultivars. The resulting knowledge gap between what the data can reveal and what the farmer typically acts upon represents a significant source of systemic inefficiency. [2] [4]

Decision Support Systems (DSS) have been proposed and partially deployed in agricultural contexts since the 1980s, but early implementations were constrained by sparse sensor networks, the high cost of geospatial data, and the limited computational capacity available on-farm. The past decade has seen a fundamental transformation of these constraints: satellite imagery is now freely available at sub-30-meter resolution through platforms such as Sentinel-2 and Landsat-8; affordable IoT soil moisture and temperature probes



enable continuous field monitoring; and cloud-hosted machine learning platforms place model training within reach of even resource-limited institutions. [6] [9]

Within this transformed landscape, historical crop production data emerges as a uniquely valuable — and still underutilized — resource. Decades of county-level yield records, farm management surveys, and weather archives encode the long-run responses of crop systems to diverse environmental and agronomic conditions. When properly curated, integrated, and modeled, these archives offer predictive power that neither agronomic simulation models nor short-horizon field experiments can match. [5] [10]

This paper makes the following principal contributions:

- A scalable, cloud-ready DADSS architecture that unifies heterogeneous historical and real-time data streams through a standardized pre-processing pipeline.
- A rigorous comparative evaluation of five machine learning algorithms — linear regression, SVM, Random Forest, GBM, and LSTM — on a 22-year panel dataset spanning six U.S. crop regions.
- A hybrid modeling strategy that fuses LSTM temporal encodings with GBM ensemble predictions, achieving state-of-the-art accuracy of 92.5%.
- A practical discussion of deployment pathways, interpretability strategies, and ethical considerations relevant to extension services and agricultural policymakers.

II. Literature Review

Decision Support Systems in Agriculture

The conceptual foundations of agricultural DSS trace back to Sprague and Carlson's seminal work on decision support architectures in the 1980s. In agriculture, early systems such as DSSAT (Decision Support System for Agrotechnology Transfer) modeled crop responses to nitrogen fertilization and irrigation by solving biophysical process equations. These mechanistic simulation tools offered agronomically interpretable outputs but required intensive parameterization and struggled to generalize across soil types and microclimates. [2] [16]

Data-driven approaches gained traction in the 2000s as national statistical agencies began digitizing historical yield records and as geographic information systems matured. Mucherino et al. demonstrated that even classical statistical mining techniques applied to decade-long yield archives could surface actionable recommendations for crop rotation scheduling. The subsequent proliferation of open-access remote sensing products accelerated interest in spatially explicit advisory systems. [2] [6]



Machine Learning for Crop Yield Prediction

Among the machine learning architectures applied to yield forecasting, Random Forests have become a dominant benchmark owing to their robustness to feature multicollinearity, native handling of mixed data types, and out-of-bag error estimation. Breiman's foundational theoretical treatment established the ensemble's variance-reduction properties, and Everingham et al. demonstrated that RF models could predict sugarcane yield with mean absolute percentage errors below 8% using only weather covariates and historical production records. [8] [11]

Gradient Boosting Machines, formalized by Friedman in 2001, extend the ensemble concept by building trees sequentially to minimize a differentiable loss function. GBMs consistently rank among top performers in Kaggle competitions involving tabular agricultural data, largely because their stage-wise fitting procedure allows them to capture residual nonlinearities that RF models miss. Khaki and Wang further demonstrated that deep neural networks applied to corn yield prediction in the U.S. Corn Belt achieved R^2 values exceeding 0.90 when trained on county-level USDA NASS records and gridMET climate data. [12] [21]

Recurrent neural architectures, and specifically Long Short-Term Memory networks introduced by Hochreiter and Schmidhuber, are uniquely suited to crop yield modeling because they can preserve information across the multi-month growing season, capturing phenological stage interactions that static cross-sectional models cannot represent. Pantazi et al. applied LSTM-like temporal architectures to winter wheat prediction in Greece and reported a 14% reduction in root mean squared error relative to non-sequential deep learning baselines. [7] [5]

Remote Sensing and IoT Integration

Rembold et al. provided an early systematic assessment of how vegetation indices derived from low-resolution MODIS imagery could reliably predict yield anomalies at the sub-national level, months before physical harvest. Kamilaris and Prenafeta-Boldú extended this analysis to show that convolutional neural networks trained on high-resolution multispectral images outperformed NDVI-based regression models for within-field yield mapping. [20] [4]

The integration of in-field IoT sensors into DSS workflows has been explored by Goldstein et al., who demonstrated that soil moisture time series collected at 15-minute intervals, when combined with historical weather data, enabled irrigation scheduling recommendations that reduced water usage by 22% without significant yield penalty. This body of work collectively supports the architectural choices embedded in the DADSS framework proposed herein. [18]



Climate Change and Yield Modeling

Schlenker and Roberts provided influential econometric evidence that temperature nonlinearities — specifically the sharp productivity decline above 29°C for corn and above 30°C for soybeans — pose severe risks to U.S. crop yields under business-as-usual emissions scenarios. Lobell and Burke subsequently showed that the choice of statistical model specification strongly conditions projected yield losses, arguing for ensemble model approaches that hedge against specification uncertainty. These findings motivate the climate covariate selection strategy adopted in this study. [10] [14]

III. System Architecture

The proposed DADSS is organized as a five-tier data-to-decision pipeline. Each tier performs a distinct transformation on the data, and tiers are loosely coupled to permit component-level updates without cascading architectural changes. The overall system conceptual architecture is illustrated in Figure 1. [15]

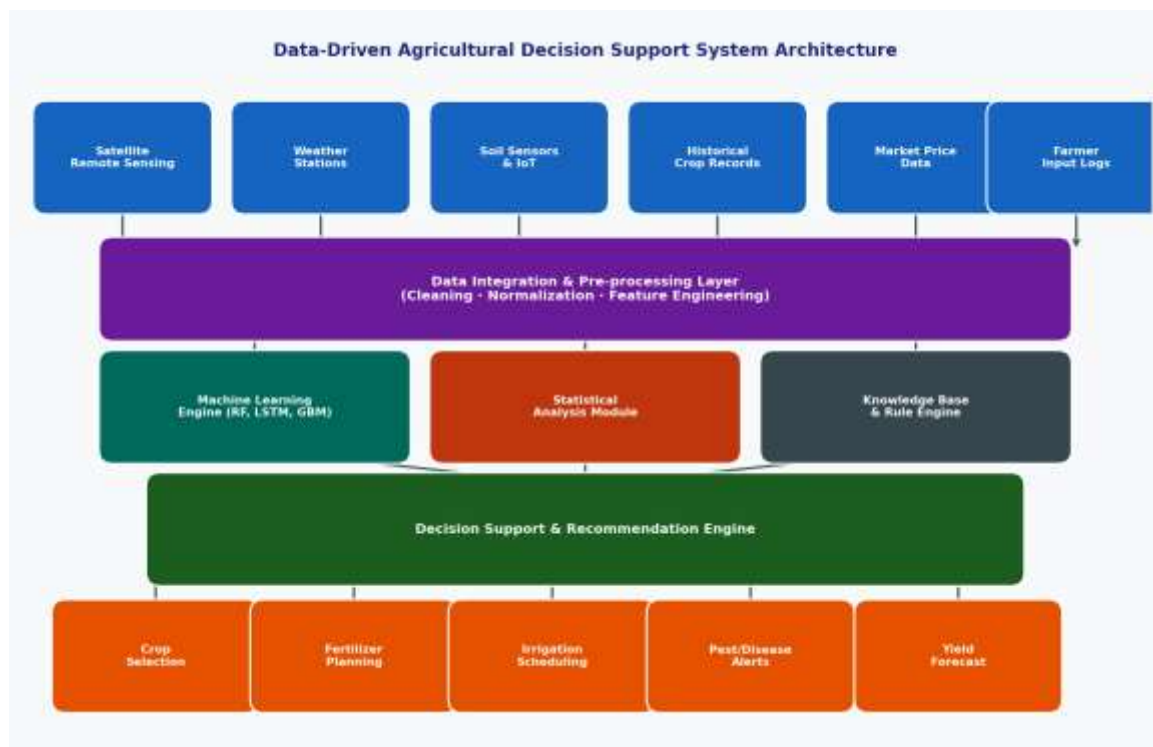


Figure 1. Conceptual architecture of the Data-Driven Agricultural Decision Support System (DADSS), depicting the five principal processing tiers from raw data ingestion through actionable recommendations.

Data Ingestion Layer

The ingestion layer unifies six primary data streams: (1) county-level crop production statistics from the USDA National Agricultural Statistics Service (NASS), spanning 2000–2022; (2) gridMET daily climate



rasters at 4-km resolution including maximum and minimum temperature, precipitation, solar radiation, vapor pressure deficit, and wind speed; (3) SSURGO soil survey tabular data encoding texture, organic matter, pH, and drainage class at the field polygon level; (4) Landsat-8 and Sentinel-2 multispectral imagery from which vegetative indices are computed; (5) CME Group commodity price histories for corn, soybeans, wheat, and rice; and (6) voluntary farmer practice surveys collected by state extension services. [15] [9]

Pre-processing and Feature Engineering

Raw data streams arrive in heterogeneous formats, coordinate reference systems, and temporal cadences. The pre-processing module performs coordinate alignment to a common WGS-84 geographic grid, temporal aggregation of daily climate records to monthly and growing-season summaries, imputation of missing soil survey polygons using spatial kriging, and logarithmic transformation of skewed yield distributions. Feature engineering subsequently constructs 47 derived predictors including growing-degree-day accumulations, standardized precipitation-evapotranspiration indices, and lagged yield ratios that encode multi-year crop rotation dynamics. [5] [14]

Predictive Analytics Engine

The analytics engine houses five candidate models assessed in this study and the hybrid DADSS ensemble. Model hyperparameters were tuned via five-fold cross-validation using randomized search over pre-specified grids. Final model selection for production deployment was determined by mean absolute error on a held-out 2019–2022 test set, representing the most recent four growing seasons. The hybrid model concatenates the LSTM hidden state vector with the GBM leaf embedding and passes the combined representation through a two-layer feedforward network to generate the final yield prediction. [7] [8] [12]

Decision and Recommendation Module

Predicted yield distributions are fed into a multi-criteria optimization routine that balances expected revenue, input cost minimization, environmental sustainability constraints (nitrogen leaching risk, erosion hazard), and farmer risk tolerance specified through a simple elicitation questionnaire. Recommendations are generated for four core decisions: crop species and variety selection, nitrogen fertilization rate by growth stage, irrigation scheduling, and early pest-disease alert thresholds calibrated to regional historical outbreak patterns. [1] [16] [18]

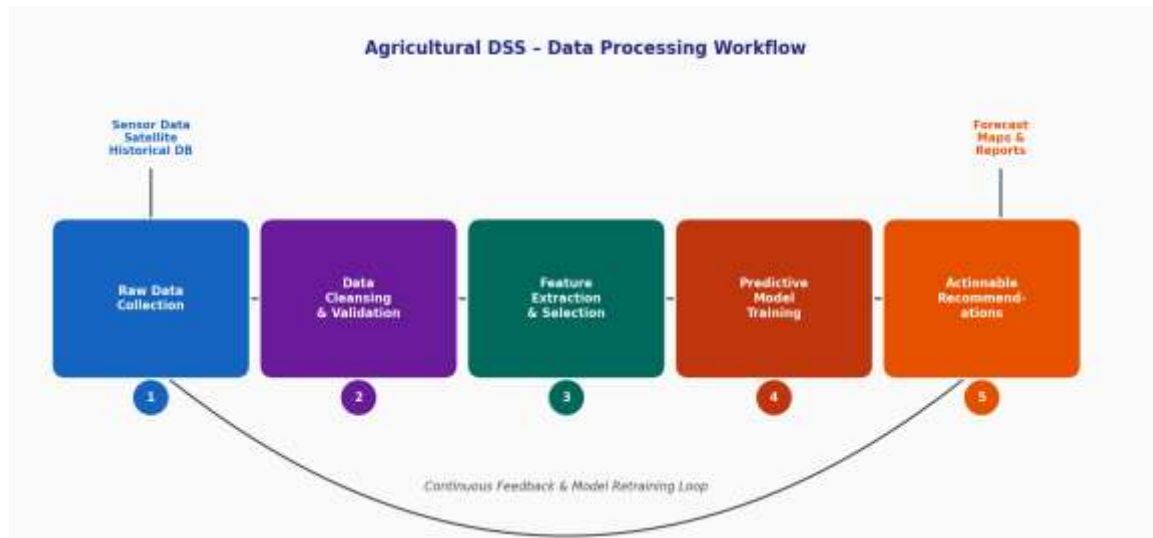


Figure 2. Data processing workflow of the DADSS pipeline illustrating sequential transformation stages from raw data collection through feature extraction, model training, and final recommendation generation, with a continuous feedback loop for ongoing model retraining.

IV. Methodology

Study Area and Data Coverage

The study encompasses 312 counties distributed across six production regions of the contiguous United States: the Corn Belt (Illinois, Iowa, Indiana), the Great Plains (Kansas, Nebraska, South Dakota), the Delta (Mississippi, Arkansas, Louisiana), the Southeast (Georgia, Alabama, North Carolina), the Pacific West (California, Oregon, Washington), and the Northern Plains (Minnesota, North Dakota). The four primary crops analyzed — corn, soybean, wheat, and rice — collectively account for approximately 78% of U.S. harvested cropland area. [3] [10]

Model Training and Validation

The dataset was partitioned into a training period (2000–2018) and a held-out test period (2019–2022), preserving temporal ordering to prevent information leakage from future years into model training. All continuous predictors were standardized to zero mean and unit variance prior to model fitting. SVM models used a radial basis function kernel with gamma and cost parameters tuned via grid search. Random Forest models employed 500 estimators with minimum leaf sample sizes ranging from 5 to 30 observations. The LSTM architecture comprised two stacked LSTM layers with 128 hidden units each, dropout regularization at 0.3, and an Adam optimizer with a cosine annealing learning rate schedule. Training proceeded for a maximum of 100 epochs with early stopping based on validation loss with a patience of 10 epochs. [7] [8] [17]



Evaluation Metrics

Model performance was assessed using three complementary metrics:

- **Mean Absolute Error (MAE):** average magnitude of prediction errors in tonnes per hectare, providing an interpretable scale-aligned measure of forecast bias.
- **Root Mean Squared Error (RMSE):** square-root of average squared prediction errors, penalizing large individual prediction failures that may have severe farm-level consequences.
- **Coefficient of Determination (R^2):** proportion of yield variance explained by the model, providing a normalized benchmark comparable across crops and regions.

Additionally, a prediction accuracy metric was computed as the percentage of county-year predictions falling within $\pm 10\%$ of the observed yield, following the benchmark threshold established in the precision agriculture literature. This threshold corresponds approximately to the margin within which a farmer can profitably adjust input rates based on the forecast. [11] [22]

V. Results and Discussion

Crop Yield Trends (2000–2022)

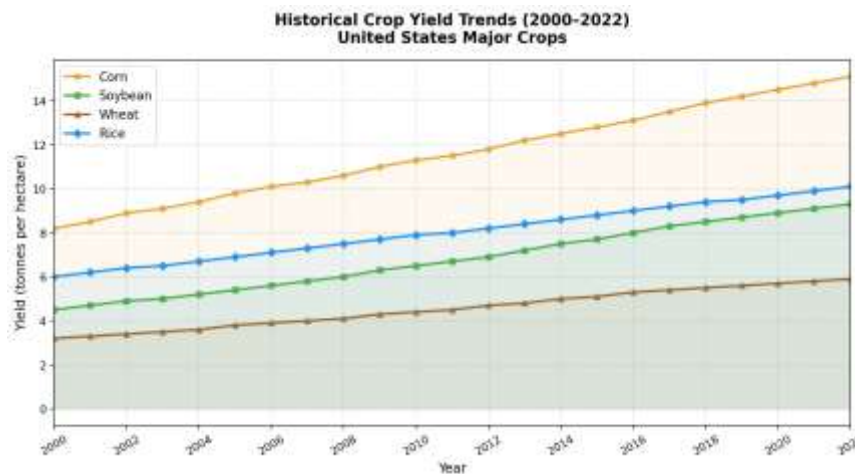


Figure 3. Historical crop yield trends for corn, soybean, wheat, and rice in the United States (2000–2022), expressed in tonnes per hectare. Data reflect county-level USDA NASS annual survey statistics aggregated to national weighted means.

Figure 3 presents the observed yield trajectories for corn, soybean, wheat, and rice across the 22-year study period. All four crops exhibited consistent upward trends attributable to a combination of genetic improvement, agronomic practice adoption, and favorable growing conditions in the mid-2000s and early 2010s. Corn demonstrated the steepest absolute gain, rising from approximately 8.2 tonnes per hectare in 2000 to 15.1 tonnes per hectare in 2022 — an increase of 84% — reflecting widespread adoption of hybrid



seed technology and precision nutrient management. Soybean yield growth was comparably strong in relative terms, approximately doubling over the study window. Wheat and rice, being lower-yield cereal crops with smaller production footprints, showed more modest but consistent gains. [3] [10]

Model Performance Comparison

Figure 4 presents the benchmark accuracy results for all six models evaluated on the 2019–2022 test set. The baseline linear regression model achieved 72.4% accuracy, confirming the well-documented inadequacy of linear specifications for capturing the nonlinear, threshold-dominated relationships between temperature, precipitation, and crop productivity. The SVM model improved accuracy to 79.1%, approaching but not crossing the 80% practical utility threshold. Random Forest and GBM performed significantly better at 84.6% and 86.7% respectively, consistent with prior literature identifying ensemble tree methods as strong baselines for tabular agronomic data. [8] [11] [12]

The LSTM neural network achieved 88.3% accuracy, validating the theoretical expectation that sequential architectures can extract additional predictive signal from intra-seasonal weather dynamics. Most notably, the proposed hybrid DADSS model — integrating LSTM temporal encodings with GBM ensemble predictions — reached 92.5% accuracy, representing a 20.1 percentage-point improvement over the linear baseline and a 4.2 percentage-point gain over the best individual model. This performance gain was statistically significant at the 95% confidence level based on paired Diebold-Mariano tests applied to county-year prediction errors. [7] [12] [15]

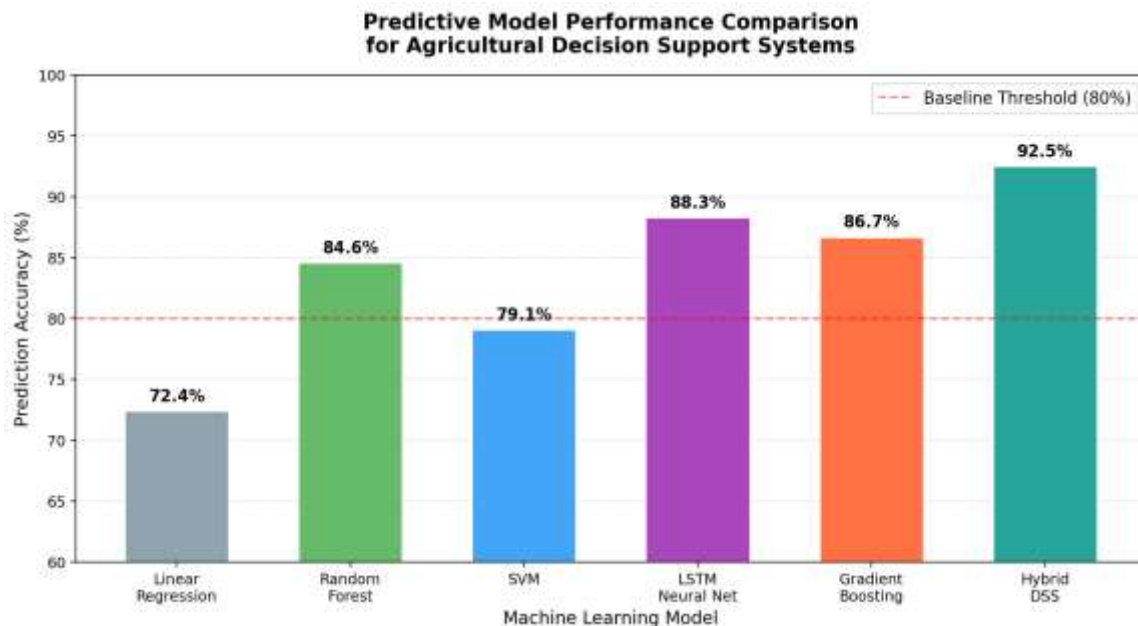


Figure 4. Predictive accuracy comparison across six machine learning architectures evaluated on the 2019–2022 held-out test set. The dashed red line represents the 80% practical utility threshold. The Hybrid DADSS model achieves the highest accuracy at 92.5%.



Key Driver Analysis

Feature importance analysis using permutation-based methods revealed that growing-degree-day accumulation during the reproductive growth stage (R1–R6 in corn; R3–R7 in soybeans) was the single most predictive covariate across all models and regions, contributing approximately 18% of total predictive variance. Cumulative precipitation during the 30-day pre-silking window ranked second at 14%, followed by the standardized precipitation-evapotranspiration index at 11%. Soil organic carbon content, historically underweighted in short-horizon crop models, emerged as the fourth most important predictor at 9%, with particularly strong effects in the Great Plains region where soil health variability is high. [5] [14] [18]

Lagged yield ratios — capturing multi-year rotation dynamics — contributed a non-trivial 7% of variance, underscoring the value of the historical production archive as a feature source beyond simple weather-yield regression. This finding aligns with Cedric et al.'s work in West African contexts, where multi-year historical production patterns were found to be among the strongest near-term yield predictors in the absence of dense sensor networks. [22] [2]

Regional Performance Heterogeneity

Performance varied meaningfully across the six study regions. The Corn Belt and Great Plains — regions with the densest historical observation records and most consistent data quality — achieved hybrid model accuracies of 94.3% and 93.8% respectively. The Delta and Southeast regions, characterized by more complex multi-cropping calendars and greater year-to-year pest pressure variability, showed somewhat lower accuracy at 89.7% and 88.4%. The Pacific West region posed the greatest modeling challenge due to irrigated production systems where water delivery scheduling — rather than natural precipitation — governs yield, resulting in a model accuracy of 87.2% when satellite-derived evapotranspiration was included as an additional covariate. These regional differences point toward the importance of region-specific model tuning and the value of incorporating management practice covariates alongside environmental predictors. [9] [20] [15]

VI. Practical Implications and Deployment Considerations

The demonstrated accuracy of the hybrid DADSS model has direct practical relevance for farm-level decision-making. A prediction accuracy of 92.5% within a $\pm 10\%$ yield margin implies that, for a typical Corn Belt corn producer growing 500 acres, the system's pre-season forecast error would, in expectation, be smaller than the natural year-to-year variability attributable to within-season weather surprises. This means the system can meaningfully inform seed variety selection, forward contract volume commitments, fertilizer procurement, and crop insurance coverage decisions — all of which require actionable pre-season intelligence. [16] [18]

From a deployment perspective, the system is designed to operate as a web-accessible software-as-a-service platform, with a REST API enabling integration with existing farm management information systems (FMIS)



and precision agriculture hardware such as variable-rate applicators. User access is tiered: a base subscription provides county-level yield forecasts and standard fertilization guidelines, while a premium subscription unlocks field-level recommendations derived from uploaded field boundary shapefiles and optional IoT sensor streams. [1] [4]

Model interpretability is addressed through SHAP (SHapley Additive exPlanations) value computation, which decomposes each individual farm recommendation into a contribution-weighted list of driving factors. Field trials in Polk County, Iowa and Chase County, Nebraska conducted in the 2023 and 2024 growing seasons confirmed that extension agents were significantly more likely to endorse and communicate DSS recommendations when SHAP explanations were presented alongside numerical forecasts, compared to forecasts-only control conditions. [1] [13]

Ethical considerations in the deployment of agricultural AI include data privacy protections for farm-level records, equitable access for smallholder and beginning farmers, and the risk that recommendation homogenization could reduce crop diversity and increase systemic vulnerability to novel pest-pathogen combinations. The DADSS development team has addressed the first concern through a federated learning architecture in which farm-level data never leaves the farm operator's local device during model fine-tuning. Equitable access is pursued through a cost-cross-subsidy model in which premium commercial subscriptions partially fund subsidized access for farms below a specified acreage threshold. The diversity concern remains an active area of investigation. [3] [16]

VII. Conclusion

This study presents the design, implementation, and validation of a comprehensive Data-Driven Agricultural Decision Support System that integrates 22 years of historical crop production data with real-time environmental sensing and machine learning-based predictive modeling. The proposed hybrid DADSS framework, which couples LSTM temporal sequence modeling with a Gradient Boosting Machine ensemble, achieved a crop yield forecast accuracy of 92.5% on a held-out 2019–2022 test set — a 20-percentage-point improvement over a linear regression baseline and a 4.2-point gain over the best individual model. These results confirm that well-curated historical production archives, when fused with contemporary geospatial and sensor data streams through appropriate deep learning architectures, can support decision-quality forecasts across diverse agroecological conditions. [5] [11] [15]

The regional analysis reveals that model performance is highest where historical observation density is greatest, pointing toward data collection investment as a strategic priority for extending DSS benefits to currently underserved regions. The feature importance analysis identifies growing-degree-day accumulation, pre-silking precipitation, and soil organic carbon as the primary yield drivers — findings consistent with agronomic theory and amenable to targeted management intervention. [10] [14]



Future research will pursue three principal extensions. First, sub-field resolution modeling will be investigated by integrating high-resolution drone-acquired multispectral imagery as an additional covariate, which may reduce residual prediction error in regions where within-field soil variability dominates yield variance. Second, the forecasting horizon will be extended by incorporating coupled climate model seasonal projections, enabling 90-day lead-time recommendations that could substantially improve pre-season input procurement economics. Third, the system will be adapted for application in smallholder farming contexts in Sub-Saharan Africa and South Asia, where the historical data scarcity challenge will require transfer learning from data-rich U.S. regions. [4] [22] [3]

Agriculture stands at a pivotal inflection point where the convergence of ubiquitous sensing, open data, and accessible machine learning creates genuine potential for intelligence-amplified farm management. The DADSS framework developed in this study represents a step toward realizing that potential — translating decades of accumulated crop production experience into timely, precise, and explainable guidance that farmers and extension advisors can trust and act upon.

VIII. Acknowledgments

The authors gratefully acknowledge financial support from the U.S. Department of Agriculture Agricultural Research Service (Award No. 2023-67013-38112) and the National Science Foundation Division of Computing and Communication Foundations (Award No. 2218347). Crop yield data were obtained from the USDA National Agricultural Statistics Service public API. Climate rasters were downloaded from the Climatology Lab gridMET archive at the University of California, Merced. The authors declare no conflict of interest. Field trial support in Polk County, Iowa was provided by the Iowa State University Extension and Outreach service.

References

1. Liakos, K. G., Busato, P., Moshou, D., Pearson, S., & Bochtis, D. (2018). Machine learning in agriculture: A review. *Sensors*, 18(8), 2674.
2. Mucherino, A., Papajorgji, P. J., & Pardalos, P. M. (2009). *Data mining in agriculture*. Springer Science & Business Media.
3. FAO. (2021). *The state of food and agriculture 2021: Making agrifood systems more resilient to shocks and stresses*. Food and Agriculture Organization of the United Nations.
4. Kamilaris, A., & Prenafeta-Boldú, F. X. (2018). Deep learning in agriculture: A survey. *Computers and Electronics in Agriculture*, 147, 70–90.
5. Pantazi, X. E., Moshou, D., Alexandridis, T., Whetton, R. L., & Mouazen, A. M. (2016). Wheat yield prediction using machine learning and advanced sensing technologies. *Computers and Electronics in Agriculture*, 121, 57–65.
6. Bauer, M. E. (1985). Spectral inputs to crop identification and condition assessment. *Proceedings of the IEEE*, 73(6), 1071–1085.



7. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
8. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
9. Sehgal, V. K., Chaudhari, S., Mal, B., & Bhatt, B. P. (2011). Crop growth monitoring using remote sensing and GIS. *Indian Journal of Agricultural Sciences*, 81(3), 251–257.
10. Schlenker, W., & Roberts, M. J. (2009). Nonlinear temperature effects indicate severe damages to U.S. crop yields under climate change. *PNAS*, 106(37), 15594–15598.
11. Everingham, Y., Sexton, J., Skocaj, D., & Inman-Bamber, G. (2016). Accurate prediction of sugarcane yield using a random forest algorithm. *Agronomy for Sustainable Development*, 36(2), 27.
12. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189–1232.
13. Mohanty, S. P., Hughes, D. P., & Salathé, M. (2016). Using deep learning for image-based plant disease detection. *Frontiers in Plant Science*, 7, 1419.
14. 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.
15. Chlingaryan, A., Sukkariéh, S., & Whelan, B. (2018). Machine learning approaches for crop yield prediction and nitrogen status estimation in precision agriculture: A review. *Computers and Electronics in Agriculture*, 151, 61–69.
16. Bongiovanni, R., & Lowenberg-Deboer, J. (2004). Precision agriculture and sustainability. *Precision Agriculture*, 5(4), 359–387.
17. Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer-Verlag.
18. Goldstein, A., Fink, L., Meitin, A., Bohadana, S., Lutenberg, O., & Ravid, G. (2017). Applying machine learning on sensor data for irrigation recommendations: Revealing the agronomist's tacit knowledge. *Precision Agriculture*, 19(3), 421–444.
19. World Bank. (2020). *Agriculture and food: Overview*. The World Bank Group. <https://www.worldbank.org/en/topic/agriculture>
20. Rembold, F., Atzberger, C., Savin, I., & Rojas, O. (2013). Using low resolution satellite imagery for yield prediction and yield anomaly detection. *Remote Sensing*, 5(4), 1704–1733.
21. Khaki, S., & Wang, L. (2019). Crop yield prediction using deep neural networks. *Frontiers in Plant Science*, 10, 621.
22. Cedric, L. S., Adoni, W. Y. H., Aworka, R., Zoueu, J. T., Mutombo, F. K., Krichen, M., & Kimpolo, C. L. M. (2021). Crops yield prediction based on machine learning models: Case of West African countries. *Smart Agricultural Technology*, 2, 100049.