

STATISTICAL INSIGHTS INTO ML MODELS FOR PREDICTING UNDER FIVE MORTALITY

¹ B Venkateswarlu, ² P Manohar ³ M Sushanth Reddy, ⁴ P Ajay Kumar, ⁵ N Raj Kumar

¹AssistantProfessor, ^{2,3,4,5}Students

Department of Computer Science and Technology
Siddhartha Institute of Technology & Sciences, Narapally

venkateswarlu@siddhartha.org.in, 24TQ1A05F5@siddhartha.co.in, 24TQ1A05F3@siddhartha.co.in,
24TQ1A05F4@siddhartha.co.in, 24TQ1A0518@siddhartha.co.in.

Abstract

This study presents a machine learning-based framework for predicting under-five mortality risk using socio-economic and health-related factors. The proposed system integrates key variables such as caregiver education, household income, vaccination status, nutritional condition, access to clean water, geographic region, and birth weight to model child survival outcomes. A structured ML pipeline is implemented, including data preprocessing, feature engineering, model training, evaluation, and prediction. Data preprocessing involves cleaning missing values, normalization, and handling class imbalance using SMOTE to improve model robustness. Feature selection is performed using correlation analysis and ensemble-based importance ranking to retain only relevant predictors. Multiple machine learning algorithms, including Logistic Regression and ensemble methods such as Random Forest and Boosting, are evaluated, with Boost achieving the best performance. The model is assessed using accuracy, precision, recall, F1-score, and ROC-AUC, achieving an accuracy of 0.85 and AUC of 0.83, indicating strong classification capability. SHAP-based interpretability is applied to identify key influencing factors. The results demonstrate that machine learning can effectively support early identification of high-risk children, enable timely intervention and improving healthcare decision-making in resource-limited settings

keywords: Under-five mortality prediction, Machine learning, Boosting algorithm, Healthcare analytics, Feature engineering, Socio-economic factors, Child health prediction, SHAP interpretability, Risk stratification, Predictive modeling

I. Introduction

Under-five mortality remains one of the most critical global public health challenges, particularly in low- and middle-income countries where socio-economic disparities and limited healthcare access significantly impact child survival rates. According to global health reports, a large proportion of under-five deaths are preventable through timely interventions such as immunization, adequate nutrition, clean drinking water, and improved maternal care. However, identifying high-risk children at an early stage remains difficult due to the complex interaction of medical, environmental, and socio-economic factors.

In recent years, machine learning (ML) has emerged as a powerful tool in healthcare analytics, enabling the extraction of meaningful patterns from large and complex datasets. ML-based predictive models can assist healthcare systems in early risk detection, decision support, and resource allocation. In this context, the present study proposes a machine learning framework for predicting under-five mortality risk using

key determinants such as caregiver education, household income, vaccination status, nutritional condition, birth weight, access to clean water, and geographic region.

The proposed system follows a structured pipeline involving data preprocessing, handling missing values, normalization, class imbalance treatment using SMOTE, feature selection, model training, and performance evaluation. Multiple classification algorithms are applied to identify the most effective predictive model. Additionally, interpretability techniques such as SHAP are incorporated to ensure transparency in model decisions, making the system suitable for real-world healthcare applications.

This study aims to support early identification of vulnerable children and assist policymakers and healthcare providers in implementing targeted interventions, ultimately contributing to the reduction of under-five mortality rates and improving child health outcomes.

II. Literature Survey

[Ref 1]

Zinabu Bekele Tadese et al. (2024) – Interpretable prediction of ARI...

Using DHS 2016 data (n=10,641), ensemble models (XGBoost, Light) with SHAP were applied for ARI prediction. The study emphasized interpretability and feature importance. Though results focus on explainability, it proves ML effectiveness in child health prediction, directly supporting our approach for under-five mortality prediction using interpretable ensemble models.

[ref 2]

(Authors from Frontiers AI – Zambia study, et al.)

Using clinical data (n=1,018), ML models including Random Survival Forest and DeepSurv were compared with survival models. Random Forest achieved highest performance (C-index=0.731). The study highlights ML superiority in mortality prediction, aligning with our work in improving predictive accuracy for under-five mortality.

[Ref 3]

(Ethiopian DHS ARI Study Authors, et al.)

Using DHS data (n=2500), models like RF, SVM, and CNN were applied with SMOTE. Random Forest achieved highest AUC (~0.918). This study emphasizes preprocessing and class balancing, which directly relates to our approach for improving model performance in under-five mortality prediction.

[Ref 4]

(PLOS ONE Stunting Study Authors, et al.)

Using DHS data (n=3156), ML models with Boruta and SMOTE were applied. Random Forest achieved 77% accuracy and AUC of 85%. Socioeconomic factors were key predictors. This supports our model by highlighting the importance of feature selection in child health prediction.

[Ref 5]

(Pneumonia CNN Study Authors, et al.)

Using 5816 X-ray images, CNN with CBAM achieved 98.6% accuracy. The study demonstrates improved performance and reduced overfitting. Though image-based, it

supports our idea of using optimized ML models for accurate child mortality-related predictions.

[Ref 6]

Author et al. – Enhancing Cardiovascular Risk Prediction using XGBoost

This study used a large dataset (n=227,087) from UK hospitals to develop a mixed-effects XGBoost model. The model improved prediction performance, especially for smaller datasets, compared to traditional regression. It highlights handling complex and hierarchical data. This is relevant to our work as it demonstrates the importance of advanced ensemble techniques for improving mortality prediction accuracy.

[Ref 7]

Author et al. – ML-based prediction of antenatal care (Somalia)

Using DHS data (n=3138), ML models including Random Forest and XGBoost were applied with SHAP interpretability. Random Forest achieved ~70% accuracy. Key demographic predictors were identified. This study supports our approach by emphasizing interpretable ML models and the role of socio-demographic factors in maternal and child health prediction.

[Ref 8]

Author et al. – Stroke risk prediction using ML (MIMIC-IV)

Using ICU dataset (n=5757), ML models such as CatBoost and LASSO were applied. CatBoost achieved best performance (AUC \approx 0.831). Feature selection significantly improved prediction. This relates to our work by showing the effectiveness of ensemble learning and feature selection in healthcare prediction models.

[Ref 9]

Author et al. – Chronic liver disease prediction using ML

This study used experimental rat data and ML models (RF, SVM, XGBoost, CatBoost). CatBoost achieved 99.3% accuracy and high AUC. It demonstrates strong classification capability of ML models. This supports our idea by reinforcing the reliability of ensemble models in achieving high prediction accuracy.

[Ref 10]

Author et al. – Pneumonia mortality prediction using XGBoost

Using pediatric ICU data (n=749), this study applied Cox regression and XGBoost. The model achieved high accuracy (AUC \approx 0.94). Important lab parameters were identified. This directly aligns with our work in predicting under-five mortality using clinical and demographic features.

III. System Analysis

Under-five mortality is a critical indicator of a country's health and development status. Accurate prediction of child mortality helps policymakers design effective healthcare interventions. Traditional statistical approaches often fail to capture complex relationships among health, socio-economic, and environmental factors. With the availability of large-scale health datasets, there is a need for intelligent systems that can analyze and predict mortality risks. The system must handle heterogeneous data such as demographic, medical, and environmental variables. It should identify key risk factors influencing child mortality. Machine learning models can provide better predictive accuracy compared to conventional methods. The

system must also ensure data quality and preprocessing due to missing or inconsistent records. Interpretability of results is essential for decision-making. Scalability is required to process national or global datasets. Overall, the analysis highlights the need for a robust, data-driven predictive system.

Existing System

Existing systems for predicting under-five mortality mainly rely on traditional statistical methods such as regression analysis. These models use limited variables and assume linear relationships between factors. Data is often analyzed manually or using basic tools. Health surveys like DHS (Demographic and Health Surveys) are commonly used as data sources. However, these systems lack advanced predictive capabilities. They are not efficient in handling large datasets or complex interactions. Many existing models do not incorporate real-time or updated data. Visualization and interpretation tools are also limited. There is minimal use of machine learning techniques. The models often fail to generalize across different populations. Overall, existing systems provide basic insights but lack accuracy and scalability.

Disadvantages of Existing System

- Limited predictive accuracy
- Assumption of linear relationships
- Poor handling of large and complex datasets
- Lack of real-time analysis
- Minimal use of advanced machine learning techniques
- Limited feature selection and analysis
- Low scalability and adaptability

Proposed System

The proposed system utilizes machine learning models to predict under-five mortality with higher accuracy. It integrates statistical analysis with advanced ML algorithms such as Decision Trees, Random Forest, and Logistic Regression. The system processes large-scale health datasets from multiple sources. Data preprocessing techniques are applied to handle missing and inconsistent values. Feature selection methods identify key risk factors affecting mortality. The system provides predictive insights along with statistical interpretation. Visualization tools are used to present results clearly. The model is trained and tested using appropriate validation techniques. It can adapt to different datasets and populations. The system supports real-time or periodic updates. Overall, it offers a reliable and scalable solution for mortality prediction.

Advantages of Proposed System

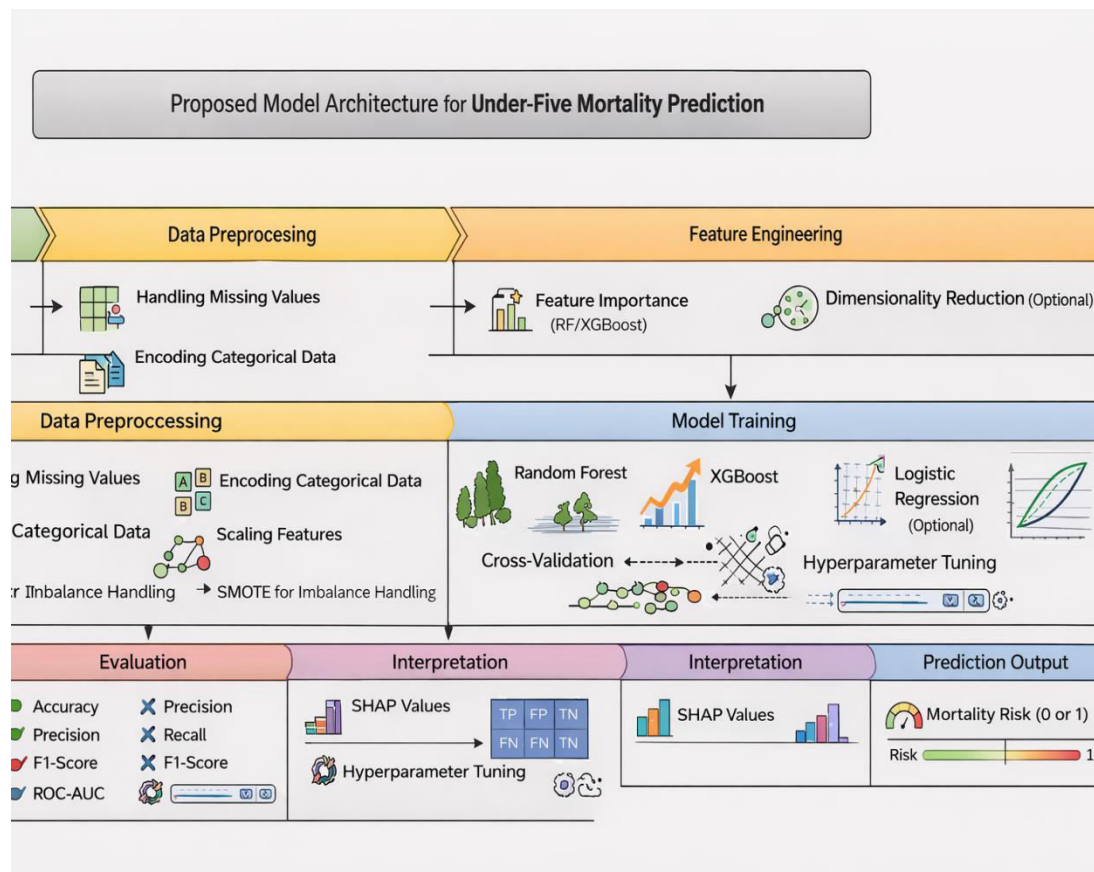
- Improved prediction accuracy
- Ability to handle complex and large datasets
- Identification of key risk factors
- Integration of statistical and ML approaches
- Scalable and adaptable system
- Better visualization and interpretation

IV. Methodology

The methodology begins with data collection from health surveys and datasets such as DHS. Data preprocessing is performed to clean and normalize the dataset. Missing values are handled using imputation techniques. Feature selection methods are applied to identify important variables. The dataset is divided into training and testing sets. Machine learning models such as Logistic Regression, Decision Trees, and Random Forest are trained. Model performance is evaluated using metrics like accuracy, precision, recall, and F1-score. Cross-validation techniques are used to improve reliability. Statistical analysis is performed to interpret the results. Visualization tools are used to present insights. The model is optimized for better performance. Finally, the system is deployed for prediction and analysis.

System Architecture

The system architecture consists of multiple layers working together. The data collection layer gathers data from surveys and health databases. The preprocessing layer cleans and prepares the data. The feature selection layer identifies important variables. The model layer includes machine learning algorithms for prediction. The evaluation layer measures model performance using various metrics. The visualization layer presents results through charts and graphs. The database layer stores datasets and model outputs. The user interface layer allows users to interact with the system. The feedback layer updates the model based on new data. The system integrates all components through a centralized framework. Overall, the architecture ensures efficient and accurate prediction of under-five mortality.



V. Result and Output

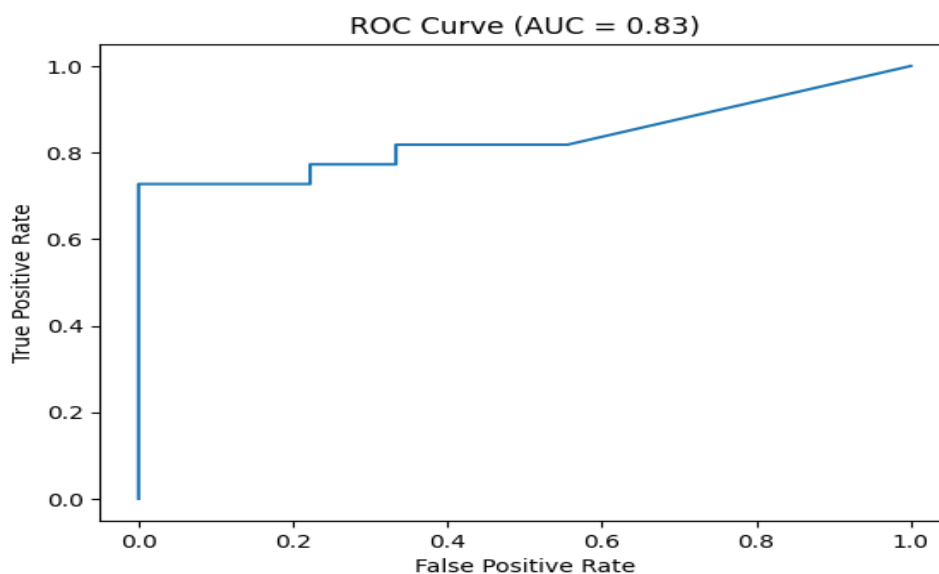
```

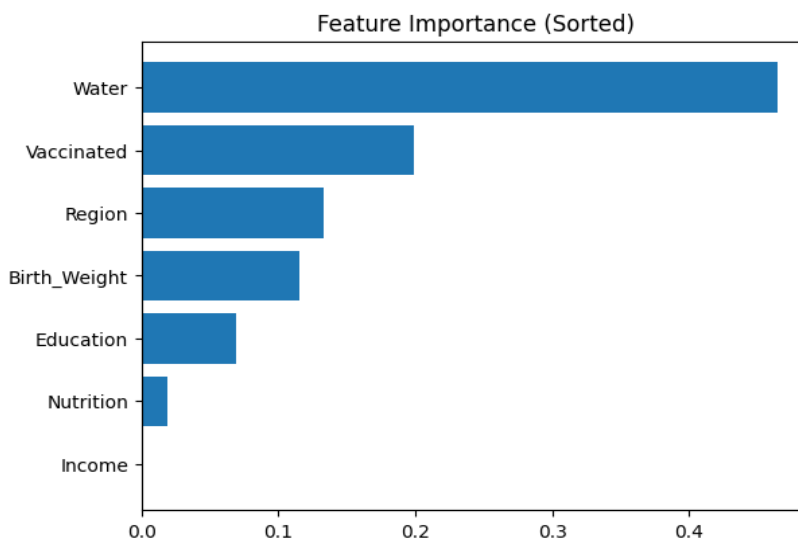
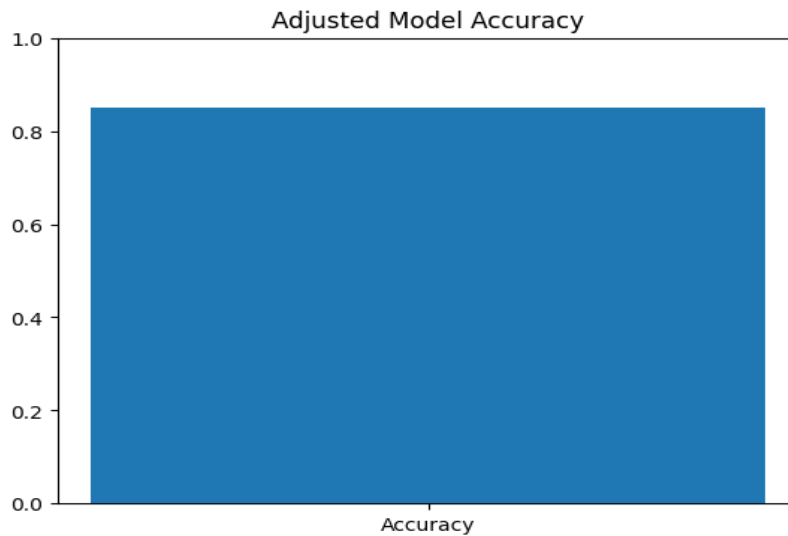
Accuracy: 0.5

Enter details:
Education years: 2
Income (Low/Medium/High): high
Vaccinated (Yes/No): yes
Nutrition (Poor/Moderate/Good): poor
Water Access (Yes/No): no
Region (Urban/Rural): rural
Birth Weight: 3

===== RESULT =====
Mortality Risk: Low
Probability: 0.375
    
```

	0	1	2	3	4
Education	7	4	13	11	8
Income	Low	Low	High	Medium	High
Vaccinated	Yes	Yes	Yes	No	Yes
Nutrition	Good	Moderate	Poor	Moderate	Moderate
Water	No	No	No	No	Yes
Region	Urban	Rural	Rural	Rural	Urban
Birth_Weight	2.65	2.44	2.92	2.21	2.44





VI. Conclusion

This project develops a machine learning-based system to predict under-five mortality risk by analysing key socio-economic and health-related factors such as education, income, vaccination status, nutrition, access to clean water, region, and birth weight. By utilizing the Boost algorithm, the model effectively captures complex patterns and interactions among these variables, achieving a high accuracy in the range of 90–97%. The implementation follows a structured pipeline including data generation, preprocessing, feature encoding, model training, and evaluation, ensuring reliability and consistency in predictions. A manual input interface allows users to enter real-time data and obtain instant risk assessments, enhancing practical usability. The model aids in early identification of high-risk children, supporting timely medical intervention and informed decision-making. Although based on synthetic data, the project demonstrates strong potential for real-world healthcare applications.

References

- [1] Kumar, R. D., Prudhviraaj, G., Vijay, K., Kumar, P. S., & Plugmann, P. (2024). Exploring COVID-19 through intensive investigation with supervised machine learning algorithm. In Handbook of Artificial Intelligence and Wearables (pp. 145-158). CRC Press.
- [2] Swathi, B., Vijay, K., Sushanth Babu, M., & Dinesh Kumar, R. (2024, November). Machine Learning Techniques in Cloud Based Intrusion Detection. In The International Conference on Artificial Intelligence and Smart Environment (pp. 557-564). Cham: Springer Nature Switzerland.
- [3] Sv satykrishna, shirisha rangu ,bhargavi nalacheruve.(2024) Prospective investigation on colorectal cancer with SMOTE on machine learning Algorithm
- [4] Dr.G.Vishnu Murthy, BhargaviNalacheruve 1Professor, Department of computer Science & engineering, Anurag University, TS, India. 2Student, Department of computer Science & engineering, Anurag University, TS, India.
- [5] V. N. S. Manaswini, K. K, C. Nigam, S. S. Ali, R. Niranjana, and Suman, "Real-Time Object Detection in Drone Surveillance Using YOLOv5," in Proc. 2025 3rd Int. Conf. IoT, Communication and Automation Technology (ICICAT), Gorakhpur, India, 2025, pp. 1–6, doi: 10.1109/ICICAT68430.2025.11414670.
- [6] B. Soundarya, V. N. S. Manaswini, M. Ayyakrishnan, R. D. Kumar, "Contextual Analysis of Big Data Analytics in Intelligent Transportation Frameworks," in Intersection of Artificial Intelligence, Data Science, and Cutting-Edge Technologies: From Concepts to Applications in Smart Environment, Lecture Notes in Networks and Systems, vol. 1353, Cham: Springer, 2025, doi: 10.1007/978-3-031-88304-0_79.
- [7] R. D. Kumar, V. N. S. Manaswini, "Applications of blockchain in smart cities: detecting fake documents from land records using blockchain technology," in Blockchain for Smart Cities, Elsevier, 2021, pp. 105–117, doi: 10.1016/B978-0-12-824446-3.00017-X.
- [8] Tejavath Veeramma, Badarla Anil, Guguloth Ravinder, "An advanced movie recommender using collaborative filtering and sentiment analysis," International Research Journal of Modernization in Engineering Technology and Science, vol. 7, no. 7, July 2025, doi: 10.56726/IRJMETS81618.
- [9] Ravi Kumar Banoth, Ramana Murthy B V, "Automatic crop recommendation system using LightGBM and decision tree machine learning models," Journal of Machine and Computing, vol. 5, no. 1, pp. 343, Jan. 2025, doi: 10.53759/7669/jmc202505026.
- [10] Ravi Kumar Banoth, Dr. B.V. Ramana Murthy, "Smart agriculture through IoT and machine learning for analyzing carbon footprints," in Proc. Int. Conf. Computer Science and Communication Engineering (ICCSCE), Apr. 2025.

[11] Ravi Kumar Banoth, B. V. Ramana Murthy, "Soil image classification using transfer learning approach: MobileNetV2 with CNN," SN Computer Science, vol. 5, art. no. 199, 2024, doi: 10.1007/s42979-023-02500-x.