# *Predicting Tanzania's Waterpoint Reliability: A Data-Driven Approach to Operational Status*



By Roy Njuguna
March 7, 2025

# SUMMARY

This project focuses on addressing water access challenges in Tanzania, where many rural and remote communities struggle to maintain functional water points. By leveraging data and machine learning, the goal is to predict the operational status of these water points, supporting efforts to improve access to safe and reliable water sources.

# BUSINESS UNDERSTANDING

Access to clean and reliable water sources is crucial for the health and well-being of communities. In Tanzania, waterpoints play an essential role in providing potable water, but many of these waterpoints face operational challenges. Some are fully functional, others operate but require repairs, and a significant number are non-functional, limiting access to safe drinking water. This project seeks to address this issue by developing a predictive model that can determine the operational status of waterpoints.

By accurately identifying which water pumps are likely to fail or need repairs, we can proactively manage resources, improve maintenance planning, and ultimately enhance access to clean water across Tanzania. This initiative will contribute to sustainable water management practices and support the long-term goal of achieving water security for all.

# DATA UNDERSTANDING

The dataset contains detailed information on waterpoints across Tanzania, capturing a variety of features that may influence the operational status of each waterpoint. These features range from geographic details to waterpoint characteristics, such as funding source, installation details, water quality, and management practices.

# MODELLING

In this project, we use XGBoost to predict the operational status of water pumps across Tanzania due to its efficiency and robustness in handling structured data. XGBoost is well-suited for classification tasks, offering strong performance and built-in regularization to prevent overfitting. To enhance the model's accuracy, we fine-tune key hyperparameters, including the number of boosting rounds (n_estimators), tree depth (max_depth), learning rate, and regularization terms (alpha and lambda). To ensure optimal performance and generalization, we apply Grid Search or Random Search with cross-validation to identify the best hyperparameter combination.

# EVALUATION

To evaluate the model's performance, we use accuracy, precision, recall, and F1-score to assess classification effectiveness. Accuracy measures overall correctness, while precision and recall provide insights into how well the model identifies each class. The F1-score balances precision and recall, making it useful for imbalanced data. Additionally, we analyze the confusion matrix to understand misclassification patterns and ensure the model performs well across all categories. By comparing these metrics across training and test sets, we assess the model's generalization and identify areas for improvement.

# RECOMMENDATIONS

Based on the model's performance and evaluation metrics, the following recommendations can enhance water point operational status predictions:

1. **Feature Engineering:** Incorporate additional relevant features, such as weather conditions, population density, and infrastructure data, to improve model accuracy.
2. **Data Quality Improvement:** Address missing or inconsistent data through imputation techniques and better data collection methods to enhance model reliability.
3. **Model Optimization:** Experiment with other ensemble methods, such as LightGBM or CatBoost, to compare performance and potentially improve results.
4. **Regular Model Updates:** Retrain the model periodically with updated data to adapt to changing conditions and maintain accuracy.
5. **Deployment and Monitoring:** Implement the model in a real-world system with continuous monitoring to assess performance and make necessary adjustments based on new data.

These steps will help refine the model, ensure better predictions, and support efforts to improve water access in Tanzania.

# NEXT STEPS

To further improve the project and ensure its practical impact, the following next steps are recommended:

1. **Expand Data Sources:** Integrate additional datasets, such as satellite imagery and real-time sensor data, to enhance model insights.
2. **Deploy the Model:** Implement the trained model in a user-friendly application or dashboard for decision-makers to monitor water point functionality.
3. **Conduct Field Validation:** Collaborate with local authorities and organizations to validate model predictions with real-world observations.
4. **Automate Data Pipeline:** Establish a continuous data collection and processing pipeline to keep the model updated with the latest information.
5. **Scale the Approach:** Explore the feasibility of applying the model to other regions facing similar water access challenges.

Thank you for reviewing this project. We appreciate your time, insights, and support in improving data-driven solutions for clean water access.

Email: roynjuguna222@gmail.com