

Tanzania Water Wells Conditions Project



Project Overview

- Tanzania, a developing country with a population of over 57 million, grapples with the substantial challenge of ensuring clean and accessible water for its people. While the country has set up numerous water points, a significant portion of these wells requires repair, and some are no longer operational. Access to clean and dependable water sources is essential for public health and economic progress.



Business Problem



The problem involves classifying well conditions into three categories:



1. Functional: Wells in good working condition, providing clean and reliable water.



2. Functional Needs Repair: Operational wells needing maintenance for sustained functionality.

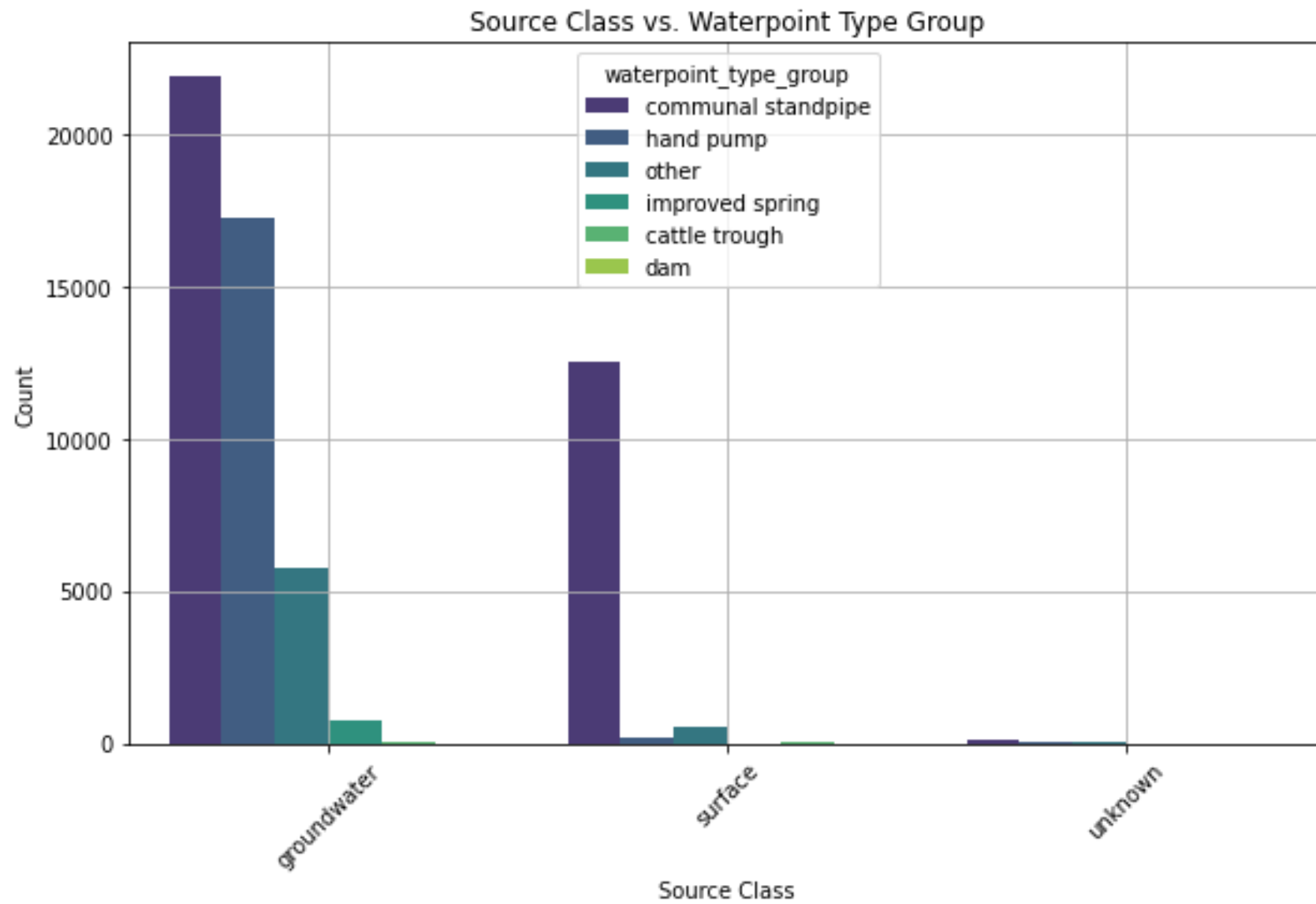


3. Non-functional: Wells that have stopped working and can't provide water.

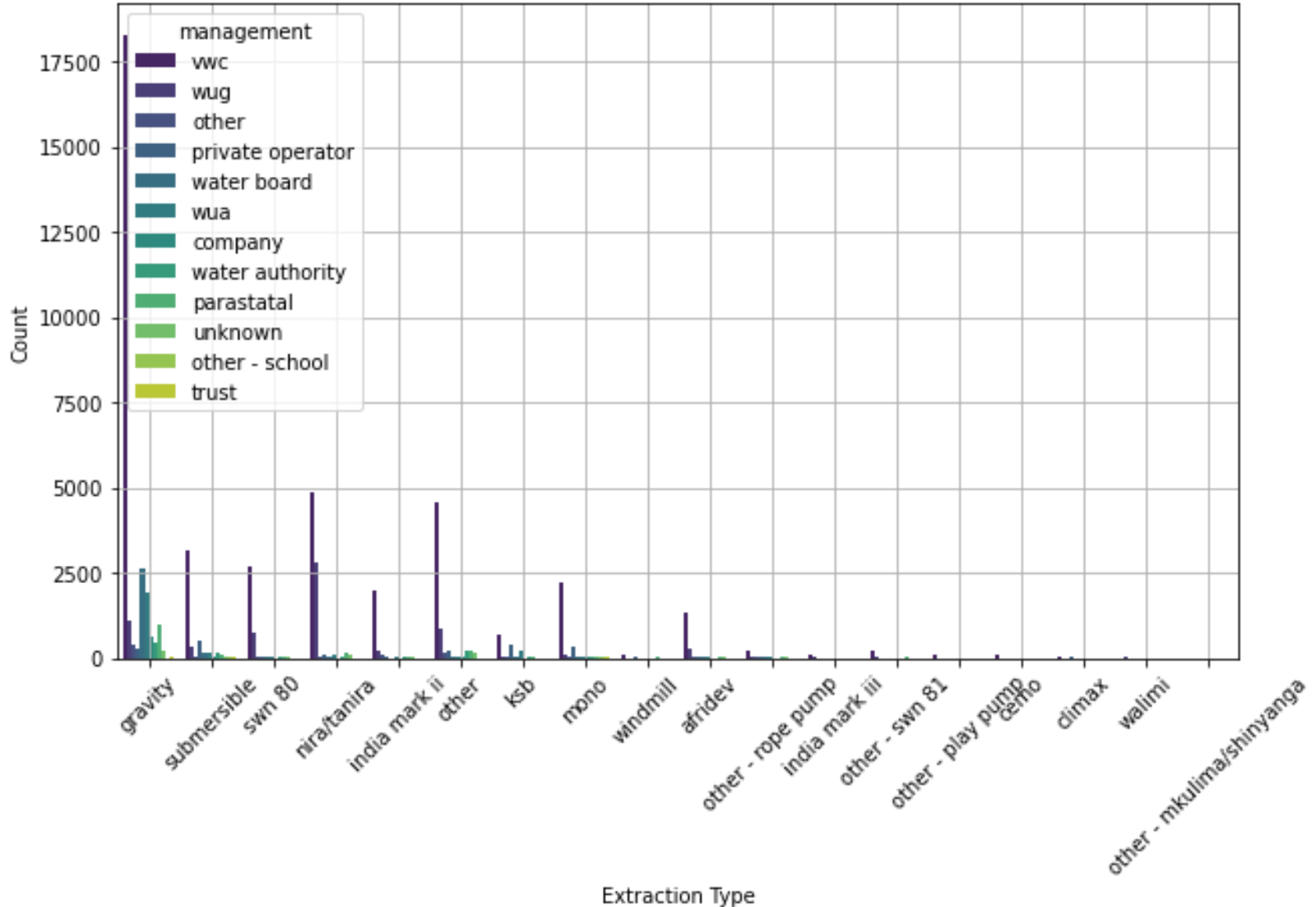
Data Understanding

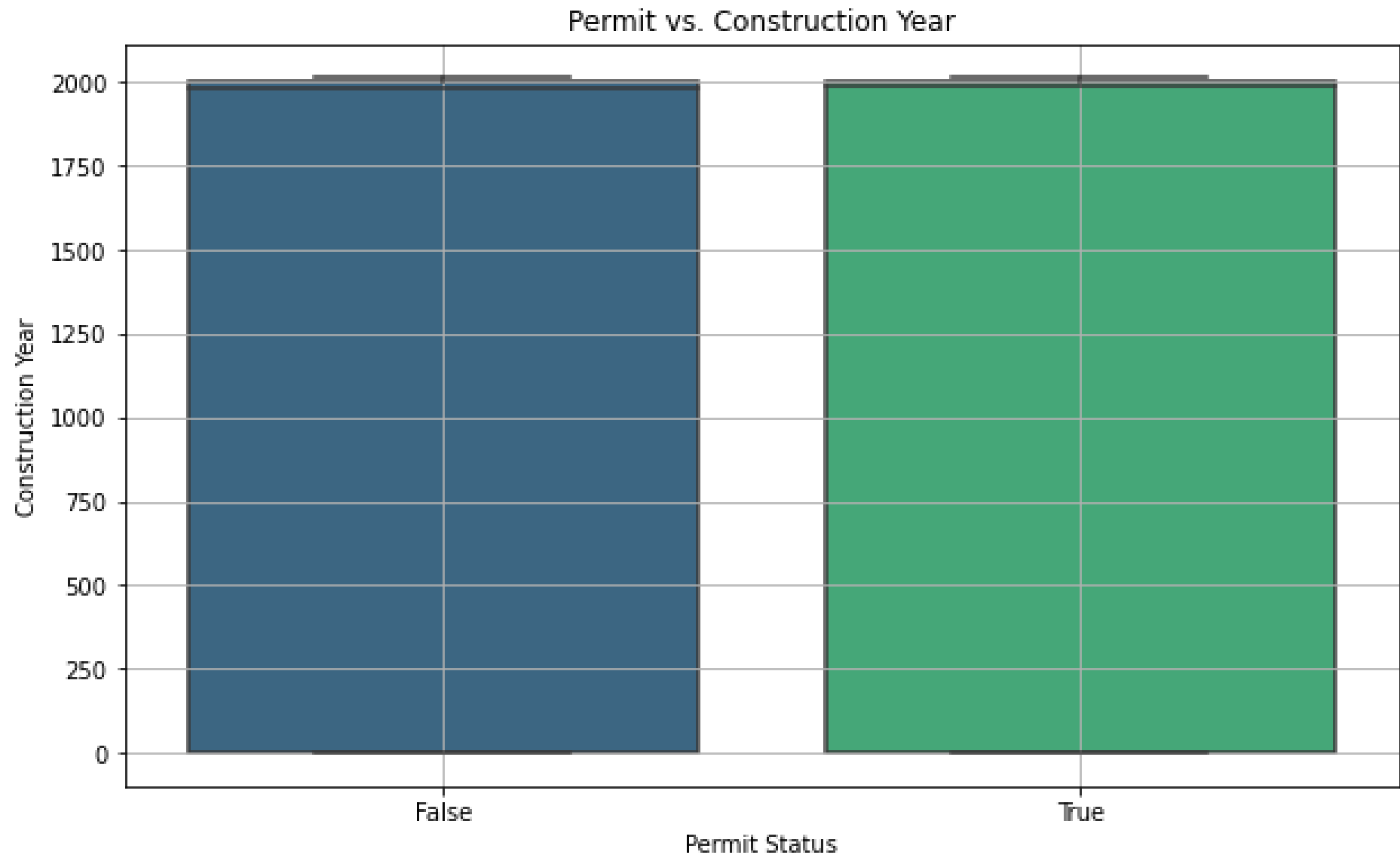
- The analysis involved a dataset from multiple sources, including Taarifa, the Tanzanian Water Ministry, and Driven Data. This dataset contains information on about 60,000 water pumps in Tanzania, with 39 attributes covering geographical, management, and technical details.
- The dataset was split into three parts:
 - 1. Test Set Values: These are the independent variables for making predictions.
 - 2. Training Set Labels: This contains the target variable "status_group" for training.
 - 3. Training Set Values: These are the features used to train the machine learning model.



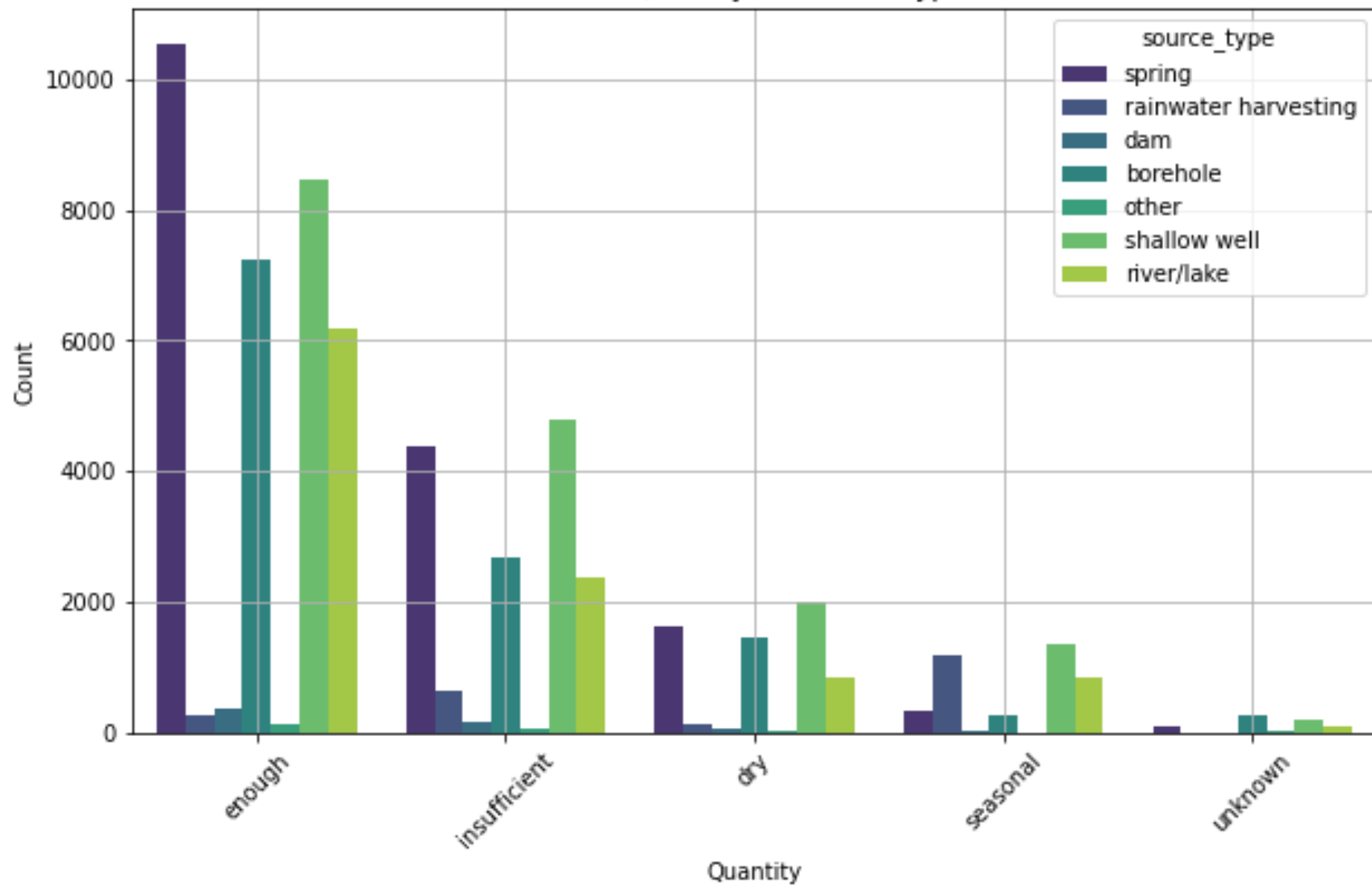


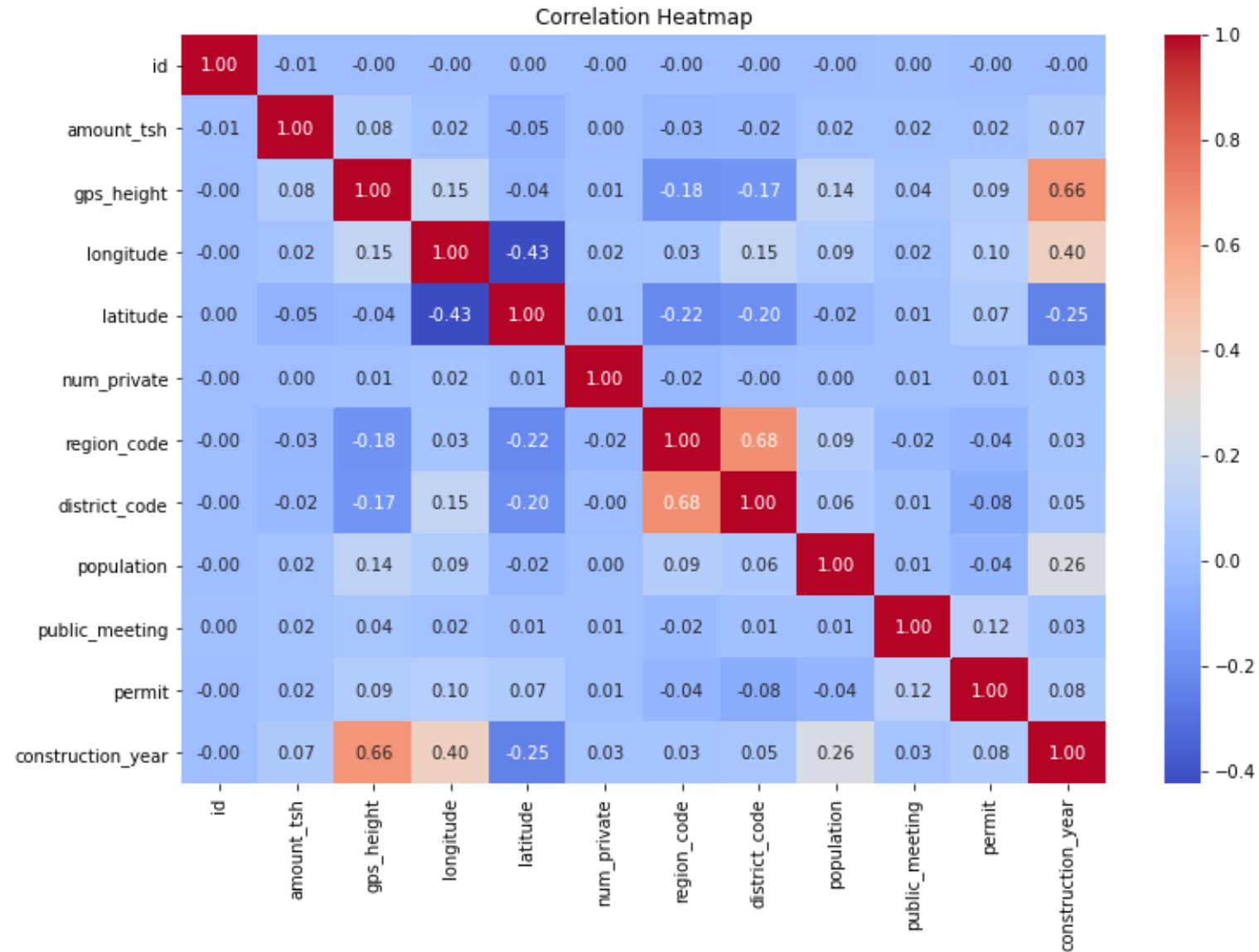
Extraction Type vs. Management





Water Quantity vs. Source Type





Modeling



AFTER TESTING VARIOUS MODELS, WE SELECTED THE ONE-HOT CODING CLASSIFIER BASED ON ITS ACCURACY SCORE. HERE'S A PERFORMANCE BREAKDOWN:



PRECISION: MEASURES ACCURATE PREDICTIONS, WITH HIGHER PRECISION INDICATING FEWER FALSE POSITIVES. THE MODEL DEMONSTRATES GOOD PRECISION FOR ALL CLASSES.



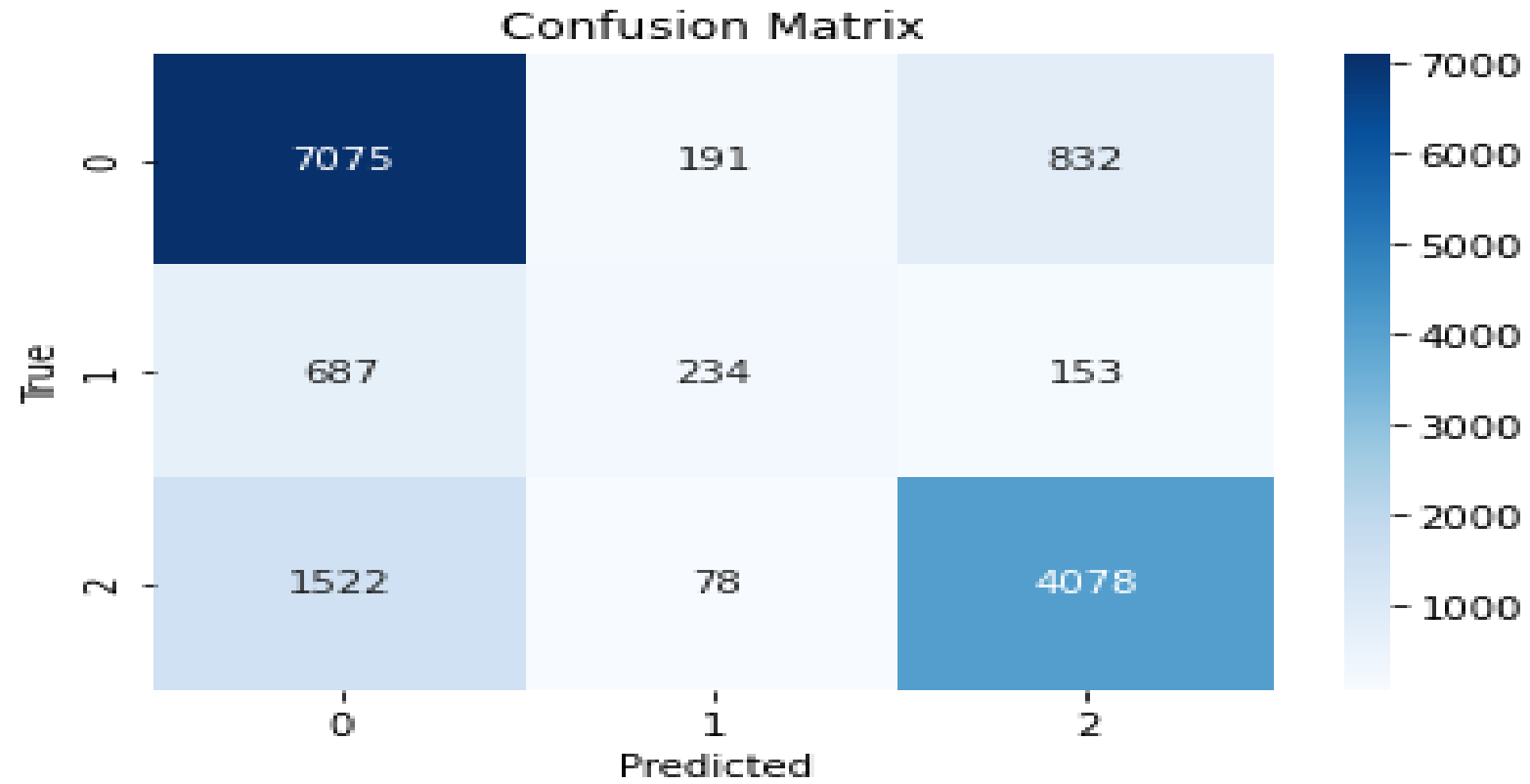
RECALL: MEASURES ACCURATE PREDICTIONS, WITH HIGHER RECALL INDICATING FEWER FALSE NEGATIVES. THE MODEL EXHIBITS REASONABLY GOOD RECALL FOR ALL CLASSES.



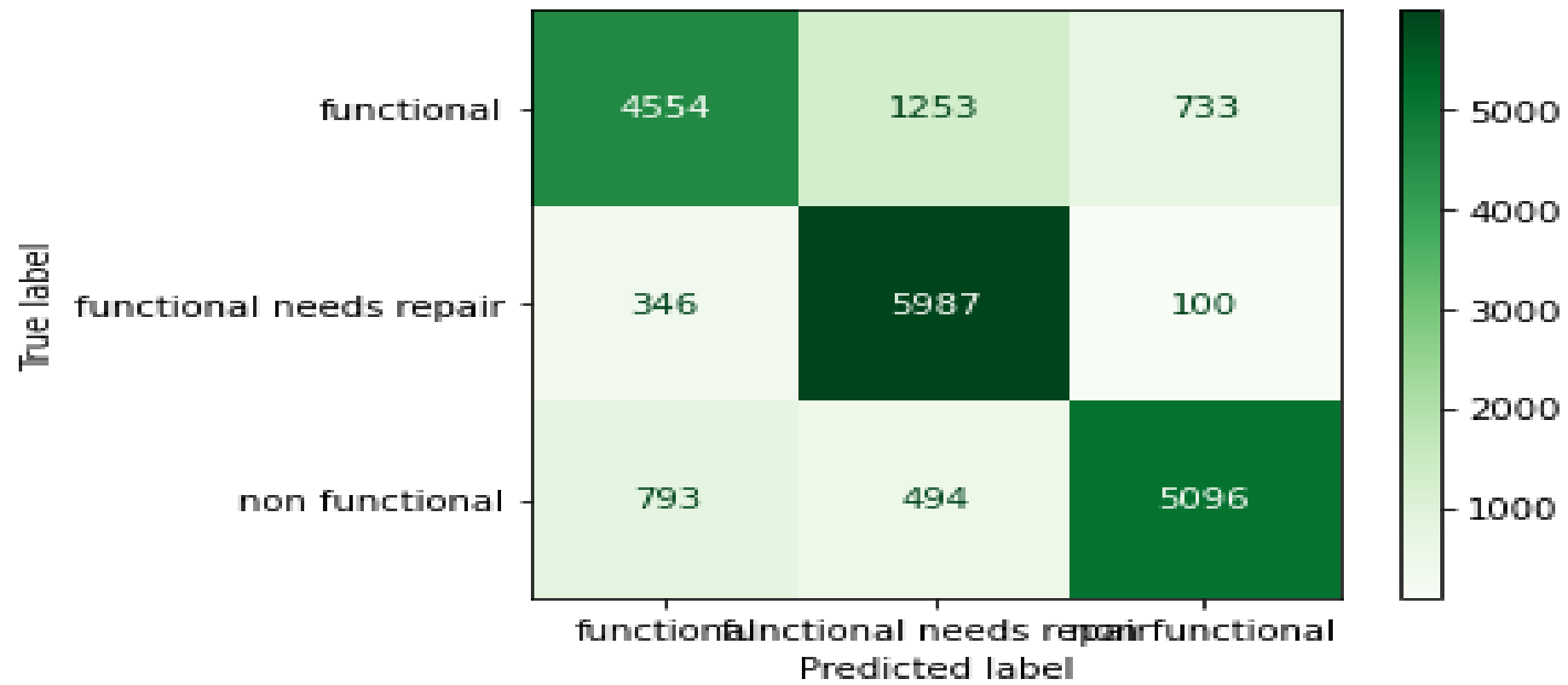
- F1-SCORE: BALANCING PRECISION AND RECALL, THE F1-SCORES ARE RELATIVELY HIGH FOR ALL CLASSES, INDICATING STRONG OVERALL PERFORMANCE.



IN SUMMARY, THE MODEL ACHIEVED AN 80.74% ACCURACY, CORRECTLY PREDICTING THE CLASS FOR ABOUT 80.79% OF INSTANCES IN THE VALIDATION SET. THE MACRO-AVERAGED F1-SCORE IS ALSO 80.79%, REFLECTING CONSISTENT PERFORMANCE ACROSS ALL CLASSES. THE WEIGHTED AVERAGE F1-SCORE IS 80.74%, CONSIDERING THE DATASET'S CLASS IMBALANCE.

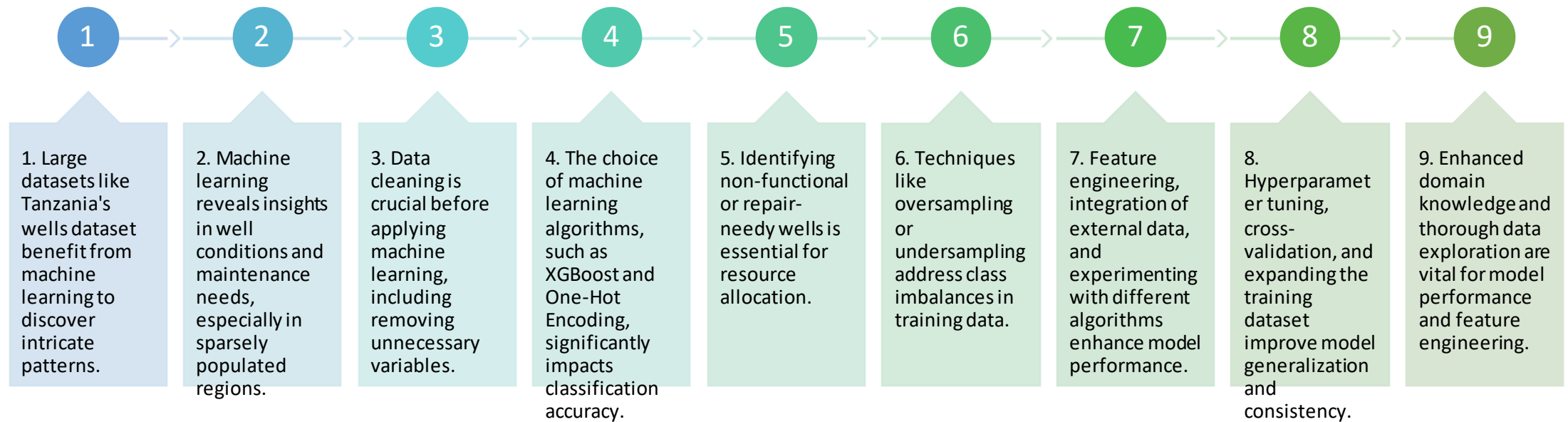


First Confusion Matrix



Second Confusion Matrix

Evaluation



Recommendations

Improving Tanzania Wells Model

- Class Imbalance Handling: Use oversampling or under sampling to address class imbalance and predict "functional needs repair" and "non-functional" wells more accurately.
- Feature Engineering: Create new features and explore external data sources to enhance predictive power.
- Algorithm Selection: Experiment with algorithms like random forests, gradient boosting, or support vector machines for better results.
- Hyperparameter Tuning: Optimize hyperparameters systematically for improved generalization and prediction accuracy.
- Cross-Validation: Implement cross-validation for robust model evaluation and reduced overfitting.
- Data Collection: Collect more data to expand the training set and improve generalization.
- Domain Knowledge & Exploration: Gain deep domain knowledge and explore data thoroughly to identify quality issues, missing values, and outliers.

Result: Improved model accuracy, particularly in identifying non-functional or repair-needy wells, aiding resource allocation for maintenance and repairs.

Next Steps

- Class Imbalance Handling: Implement oversampling or under sampling to address class imbalance and enhance predictions for "functional needs repair" and "non-functional" wells.
- Feature Engineering: Invest in creating new features and exploring external data sources to boost predictive power.
- Algorithm Selection: Experiment with algorithms like random forests, gradient boosting, or support vector machines for improved model performance.
- Hyperparameter Tuning: Optimize hyperparameters systematically for better generalization and prediction accuracy.
- Cross-Validation: Implement cross-validation for robust model evaluation and to prevent overfitting.
- Data Collection: If possible, collect additional data to expand the training set and improve model generalization.
- Domain Knowledge & Exploration: Deepen domain understanding and conduct thorough data exploration to identify quality issues and feature-target relationships.

Implementing these steps will result in an enhanced model, particularly in identifying non-functional or repair-needy water pumps, aiding resource allocation for maintenance and repairs.

Thank You

