# WATER IS FOR EVERYBODY

Bringing clean water to Tanzania.
Presented by Shefat Moral, Gabriel Santorelli, David Jimenez

# Providing functional waterwells❓❓

Our goal for this project was to build an effective ML model that can predict water pump functionality for the purpose of bringing clean water to the people of Tanzania .

EDA Process:

Analyze the data.

Select determining factors for our model.

Cleaning the data for our model.

Model Selection:

Preprocessing for our baseline model.

Decision Tree - baseline model.

A Model the People Can Trust:
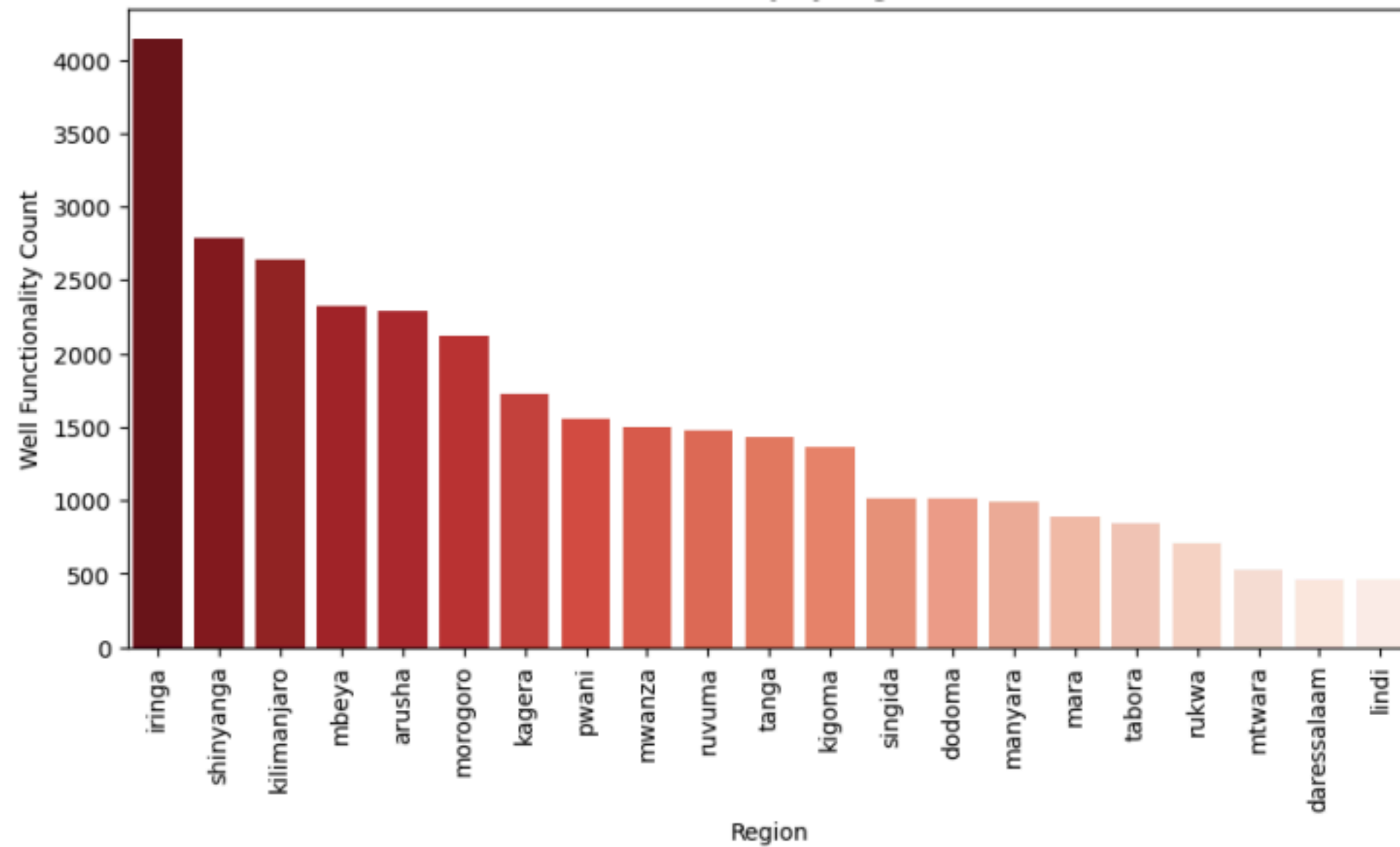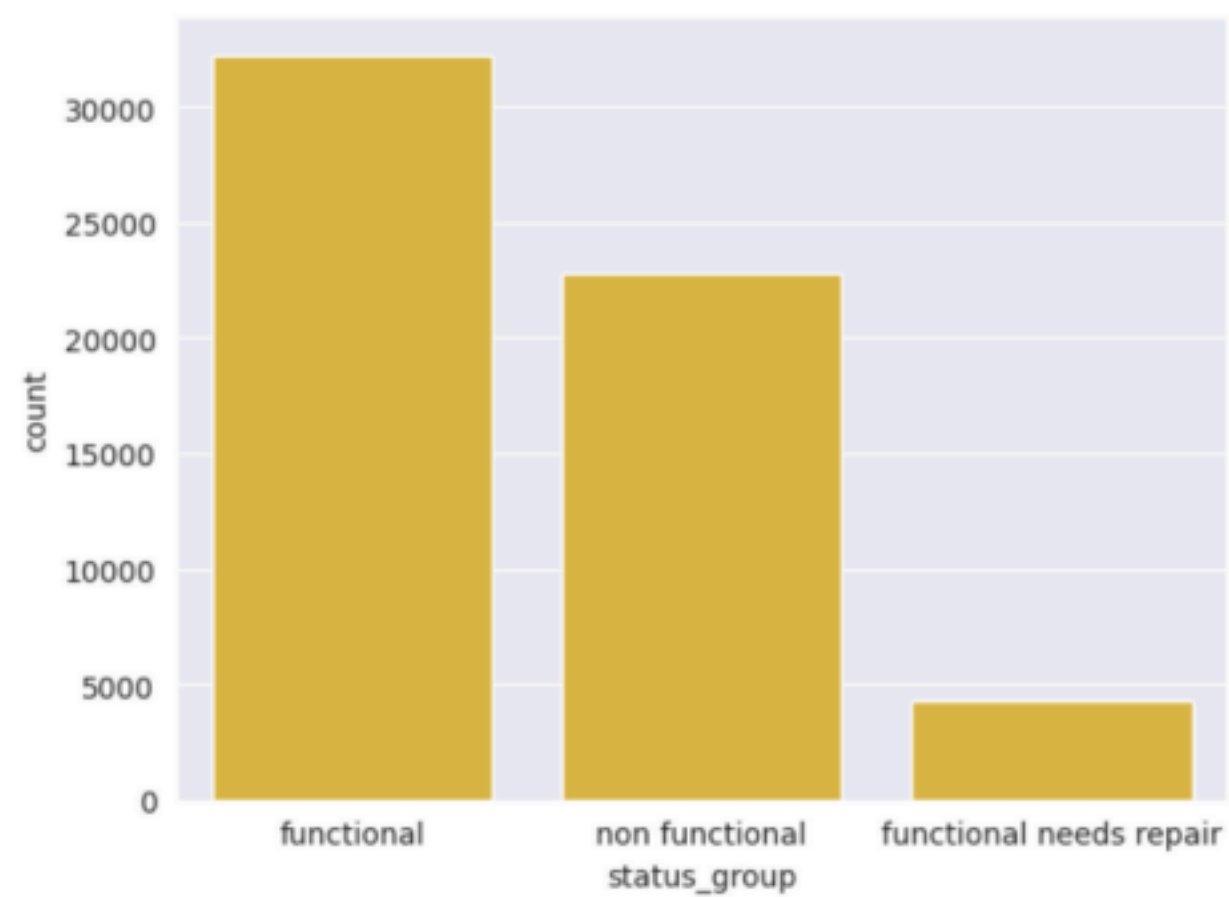
Feature Engineered

Hyperparameter Tuning

Class Balancing
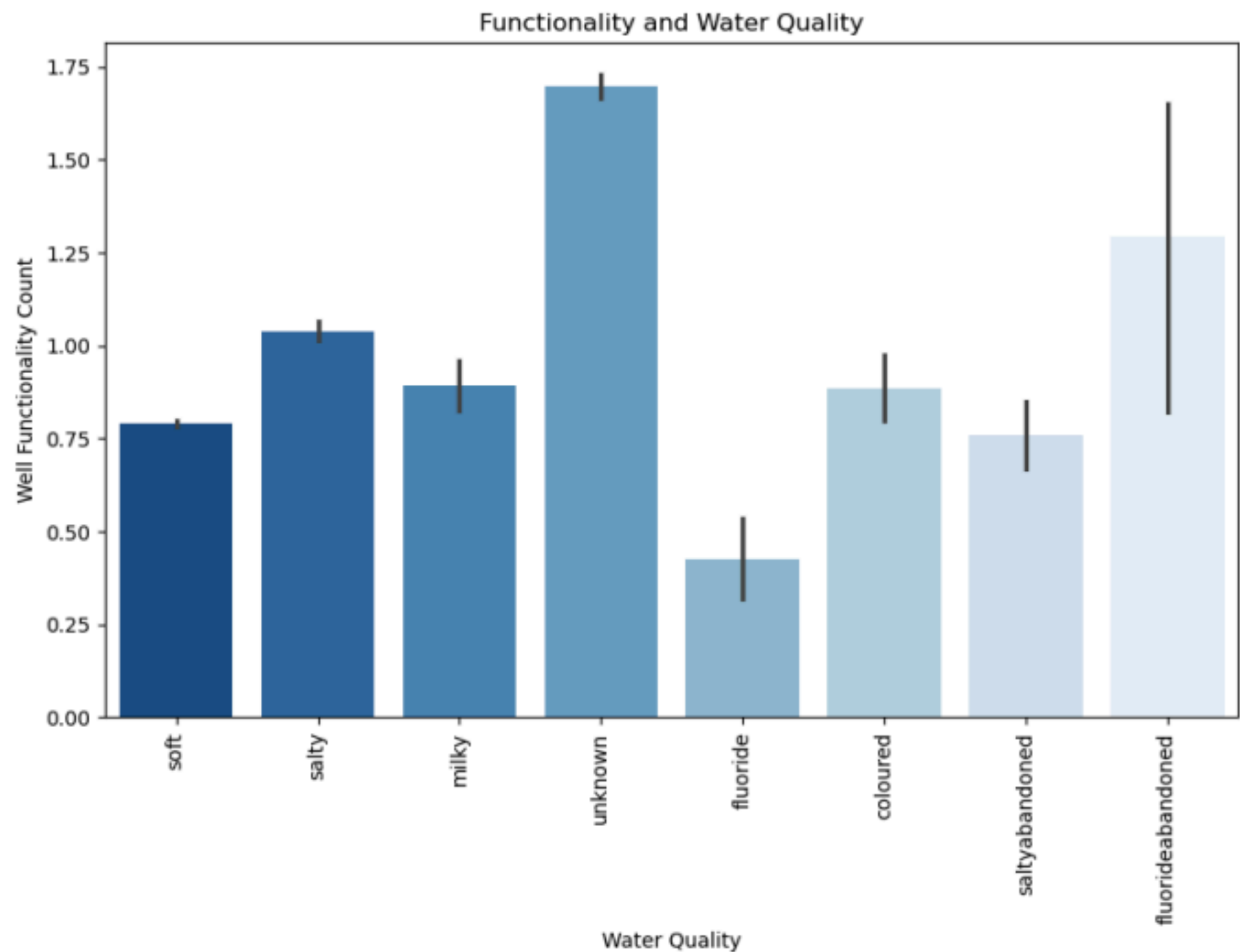
# EDA Process
## Data Exploration

Multiclass data:

    Functional

    Not Functional

    Functional Needs Repair

Looked at features like Functionality
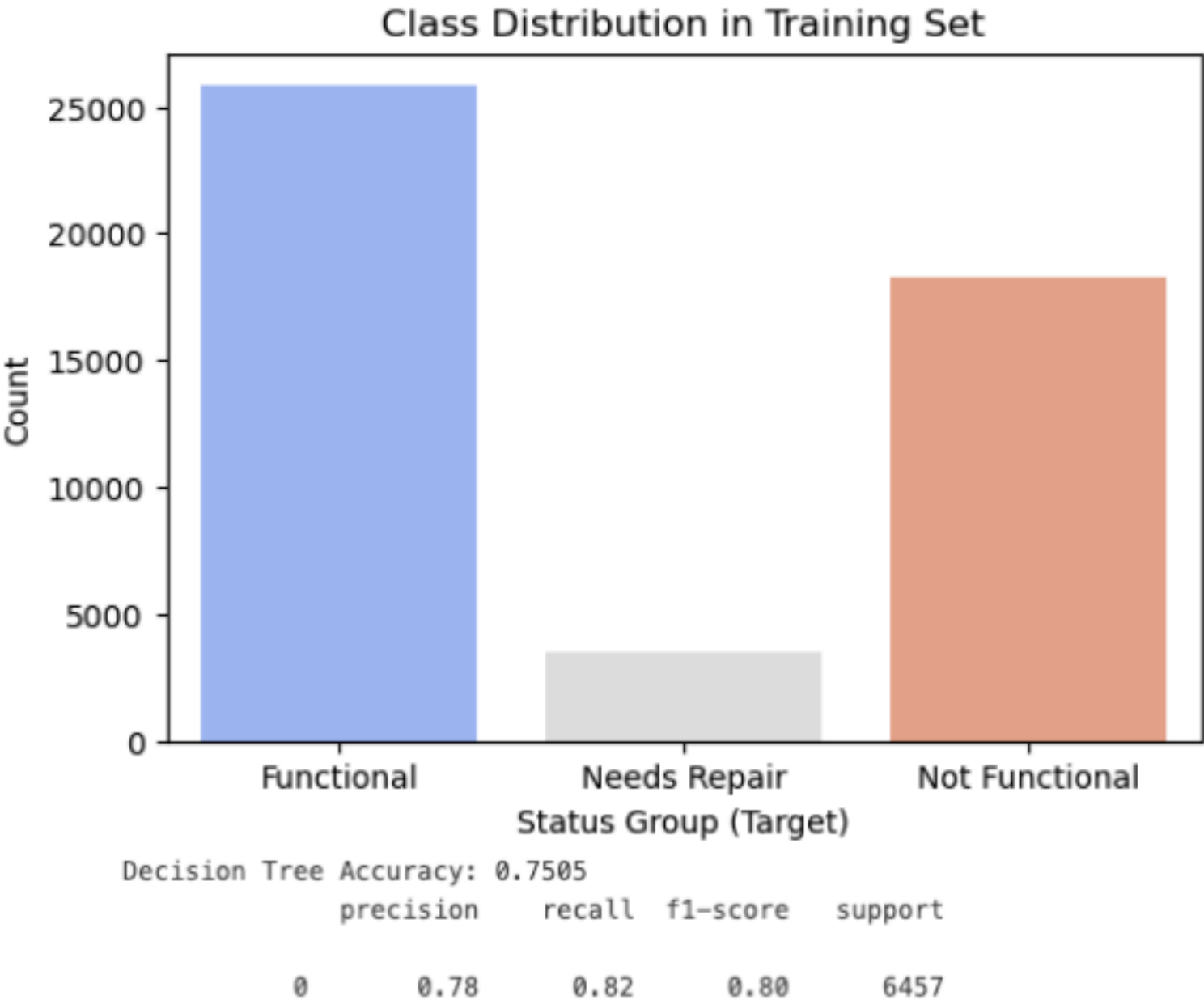
by Region

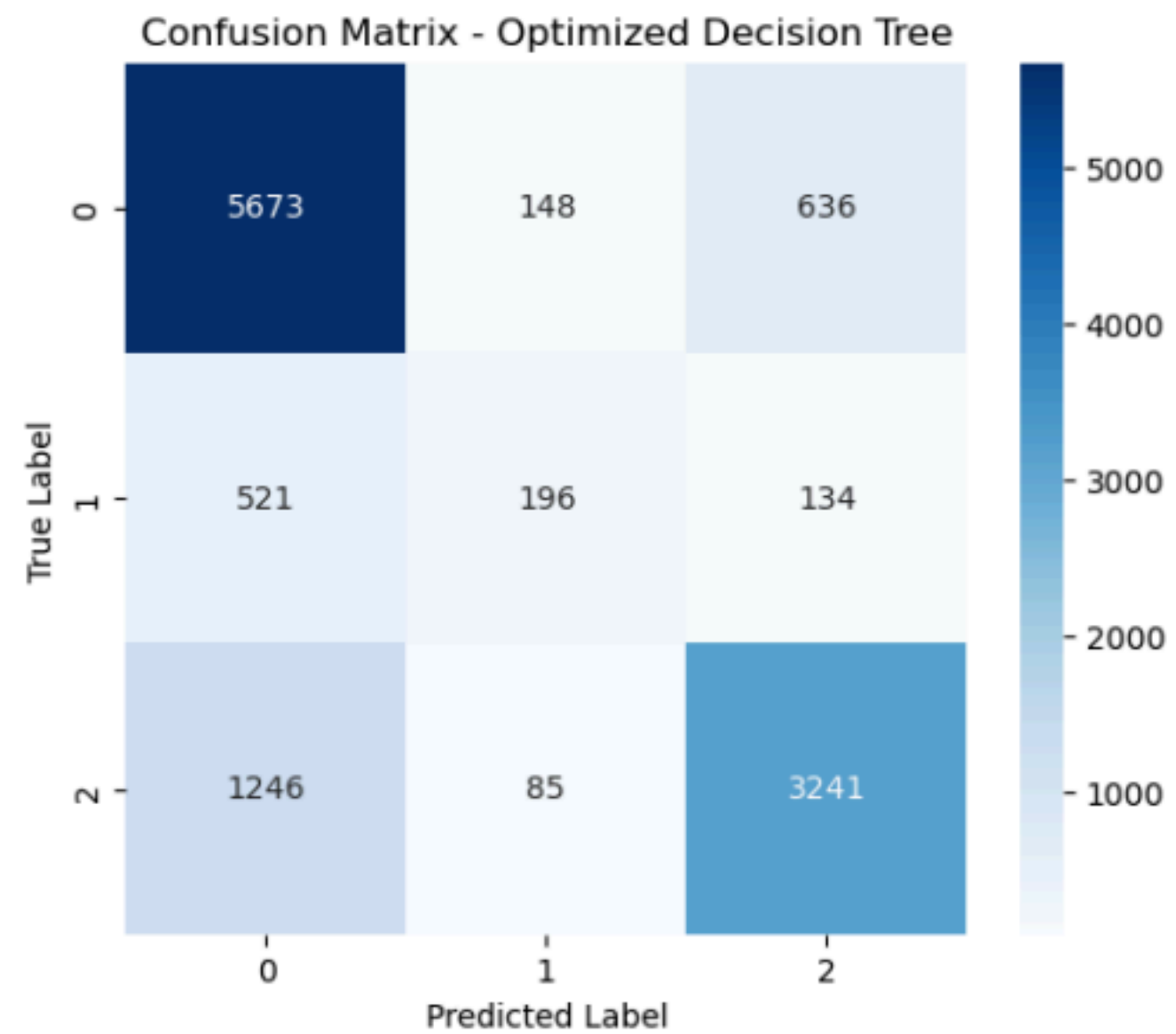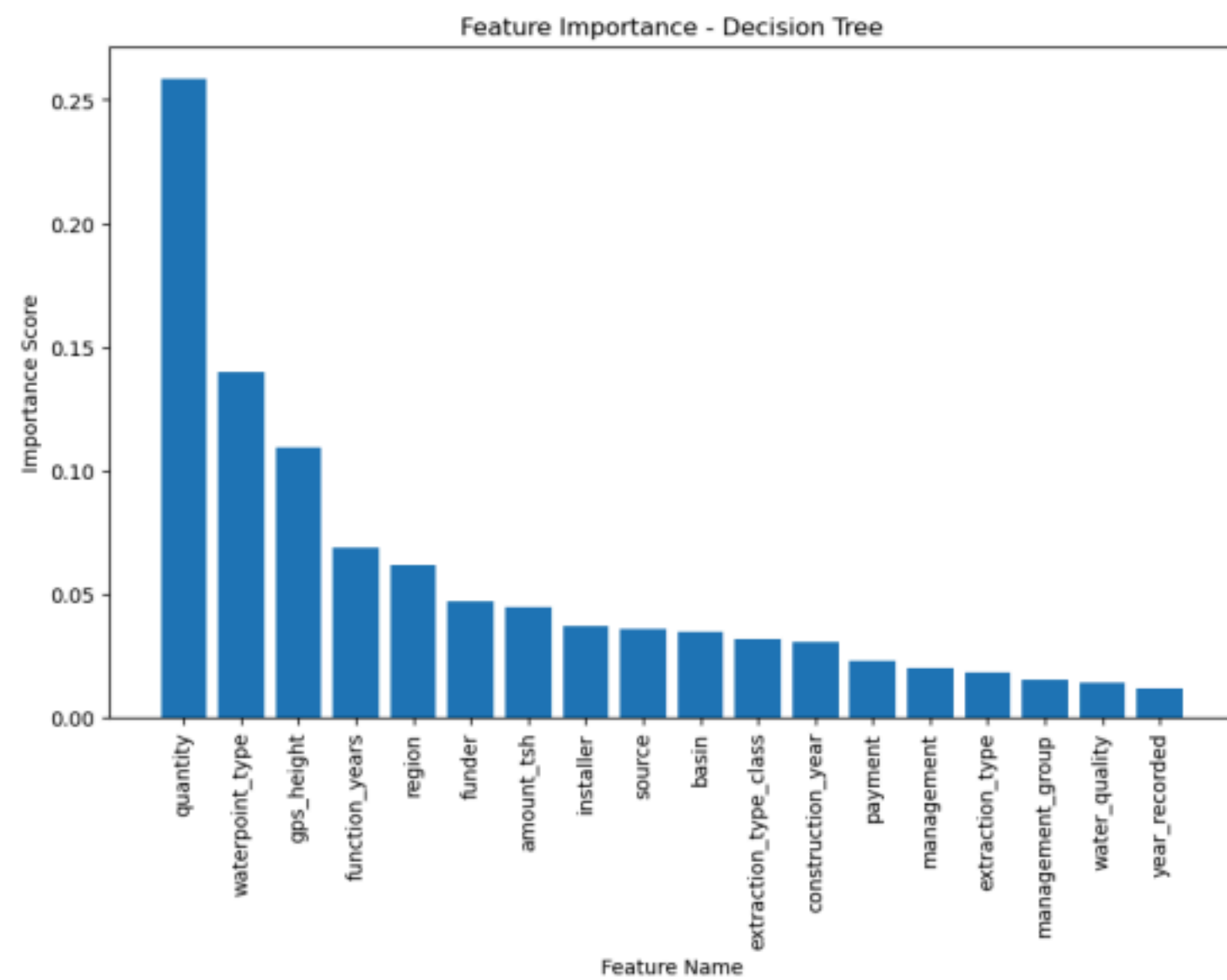Tanzania Waterwells Project

# Water Qualtiy Analysis

## Functioning waterwells + Water Quality

Functioning did not always mean clean water.
Unknown, Flourdide Abandoned, and Salty top
the Water Quality list.
  Soft(good) was near the bottom of the
quality list.



Functionality and Water Quality

Tanzanian Waterwells Project

## Model Selection



Class Distribution in Training Set

Decision Tree Accuracy: 0.7505

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.78 | 0.82 | 0.80 | 6457 |

Decision Tree Model



Feature Importance - Decision Tree
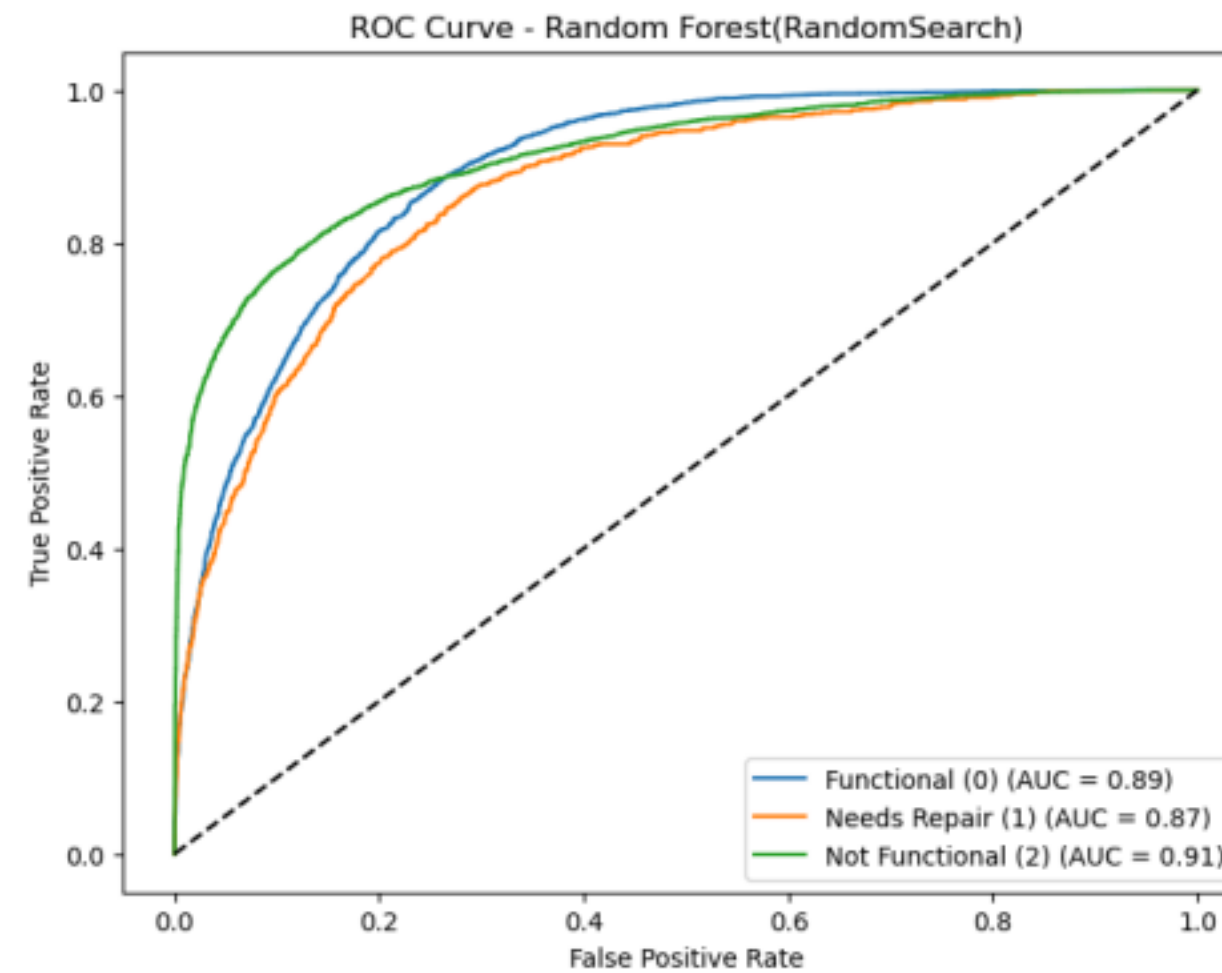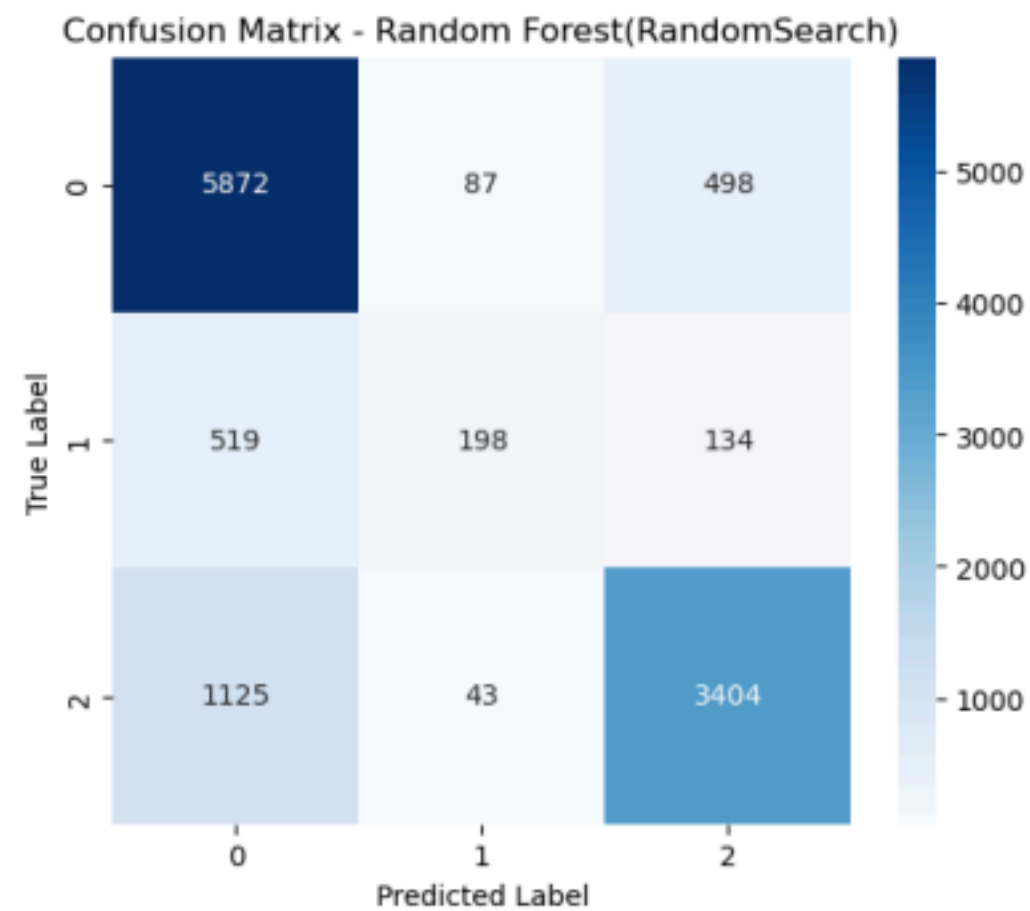
Confusion Matrix - Optimized Decision Tree

Max accuracy score of 76% 23% recall and 46% precision  False positives increased

Further EDA showed model's feature importance Tuned hyperparameters    Weighted the model, and used SMOTE

Tanzania Waterwells Project



Random Forest was the optimal model
Initial accuracy score was 79%
Increase in precision score for Class 1 (60%) and Class 2 (84%)
Recall score was 23% for Class 1
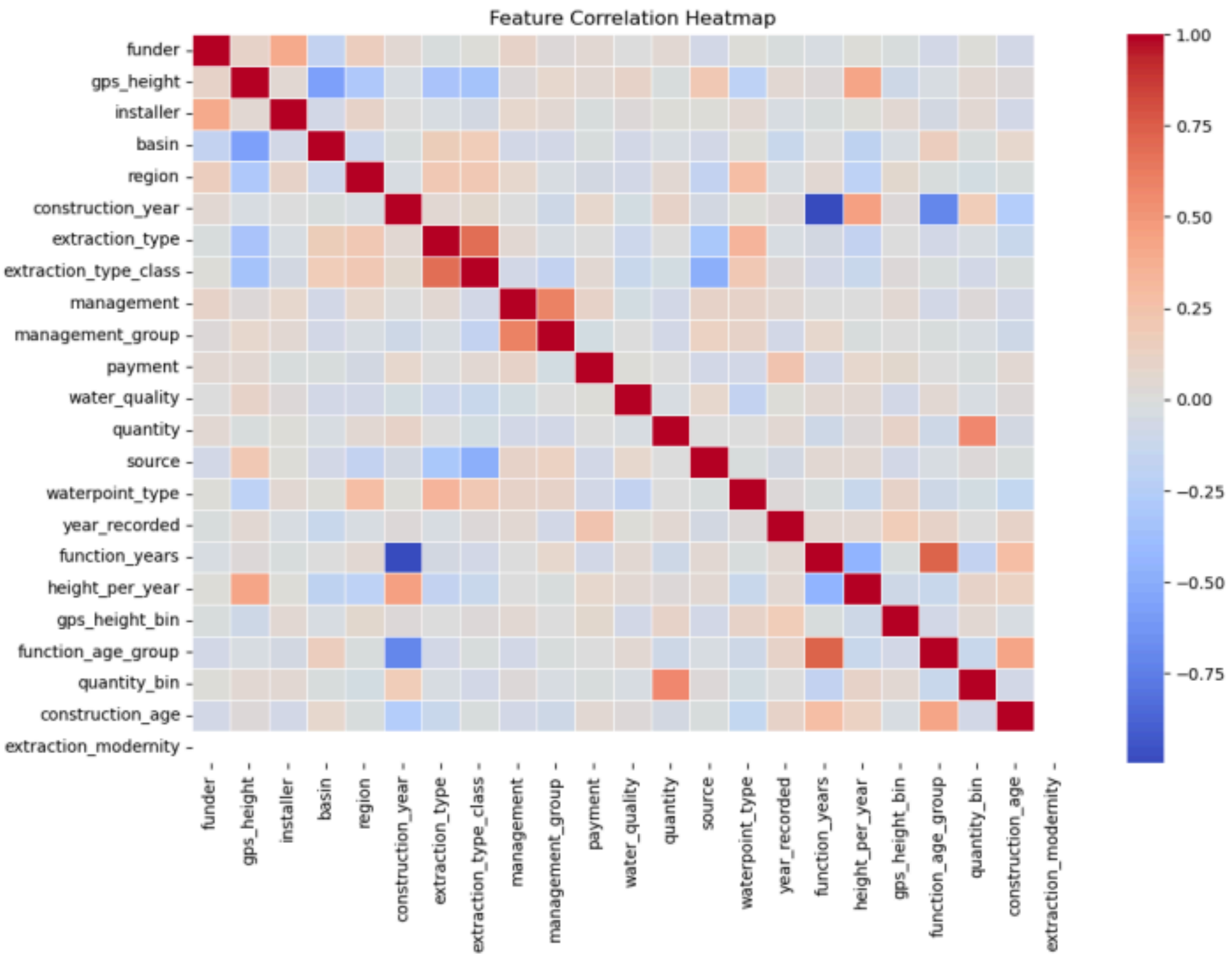
Tanzania Waterwells Project

Model identified Class 0 and Class 2 very well
Model struggled to distinguish Class 1 from other classes.

# Correlating Features
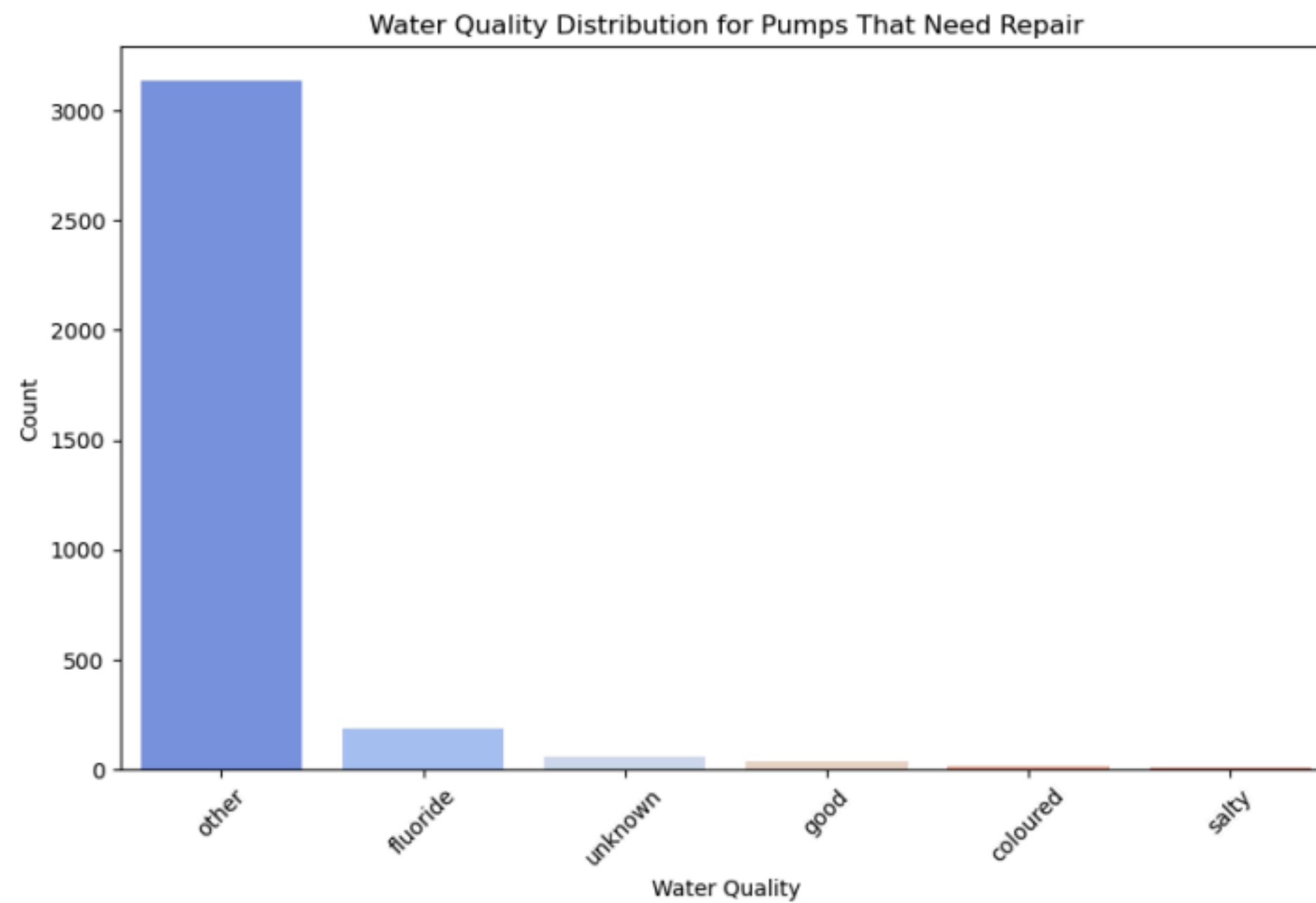
## Identify Correlating Features

Engineer features based on correlation
Binned features like gps_height and
function_years.
Raised accuracy score to 80%
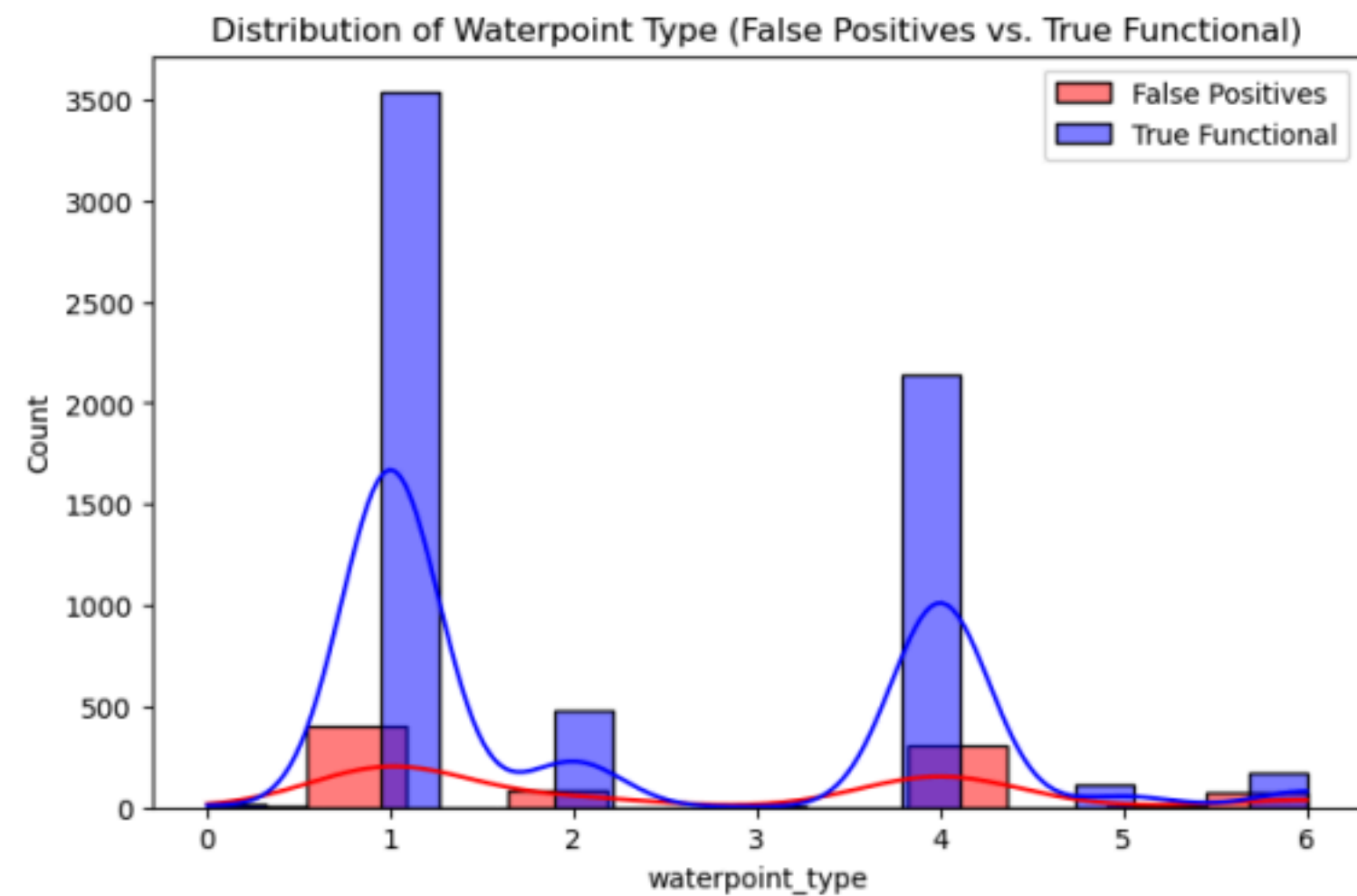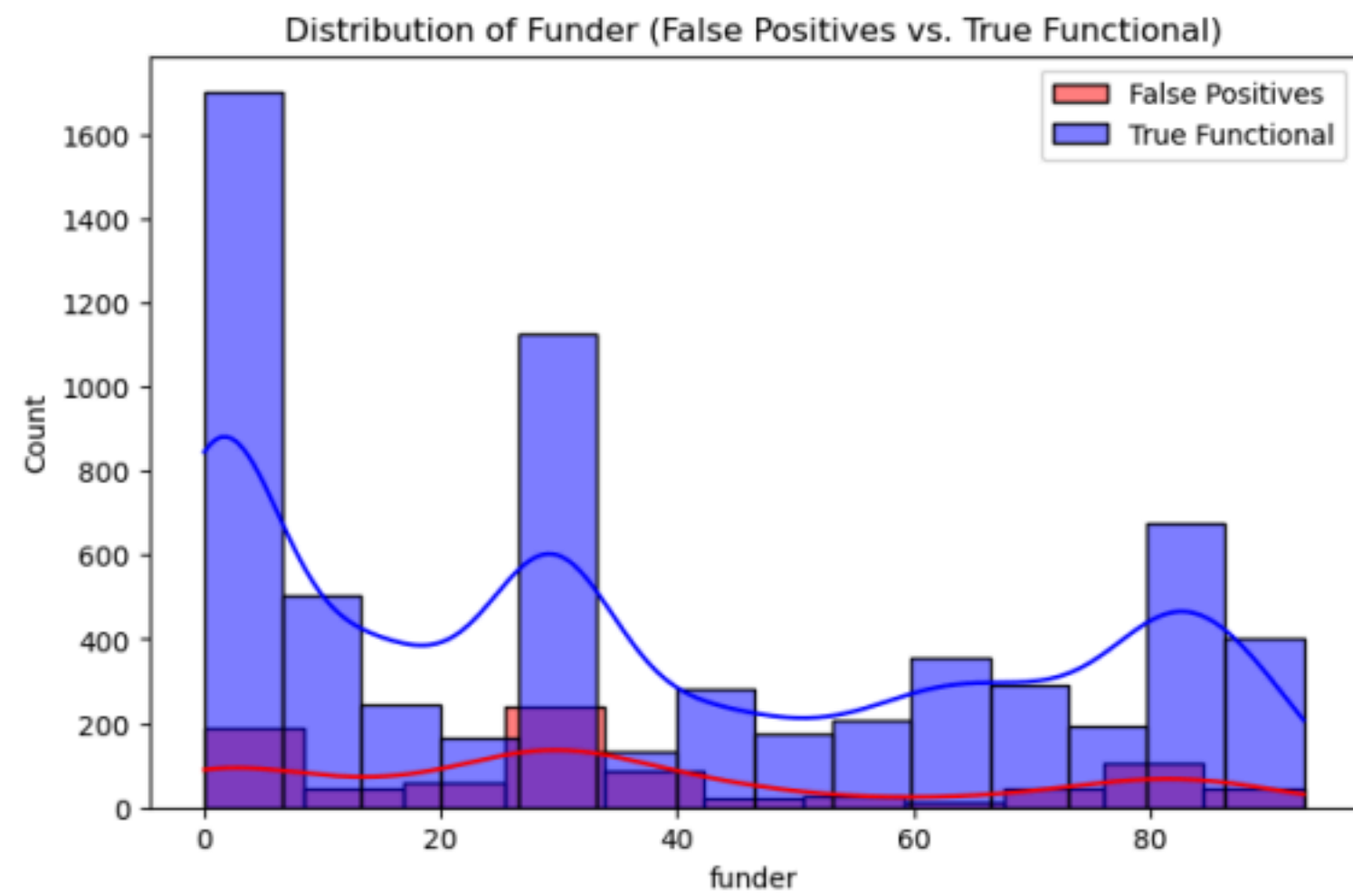
# But we could do better...

Tanzanian Waterwells Project

# Random Forest Error Analysis

Water Quality Distribution for Pumps That Need Repair

Dived into false positives

Broken pumps = poor water quality

Grouped Class 1 with Class 2

Model's accuracy score was 81%

200 less false positives

.0025 STD shows model is stable

Tanzanian Waterwells Project

Distribution of Funder (False Positives vs. True Functional)


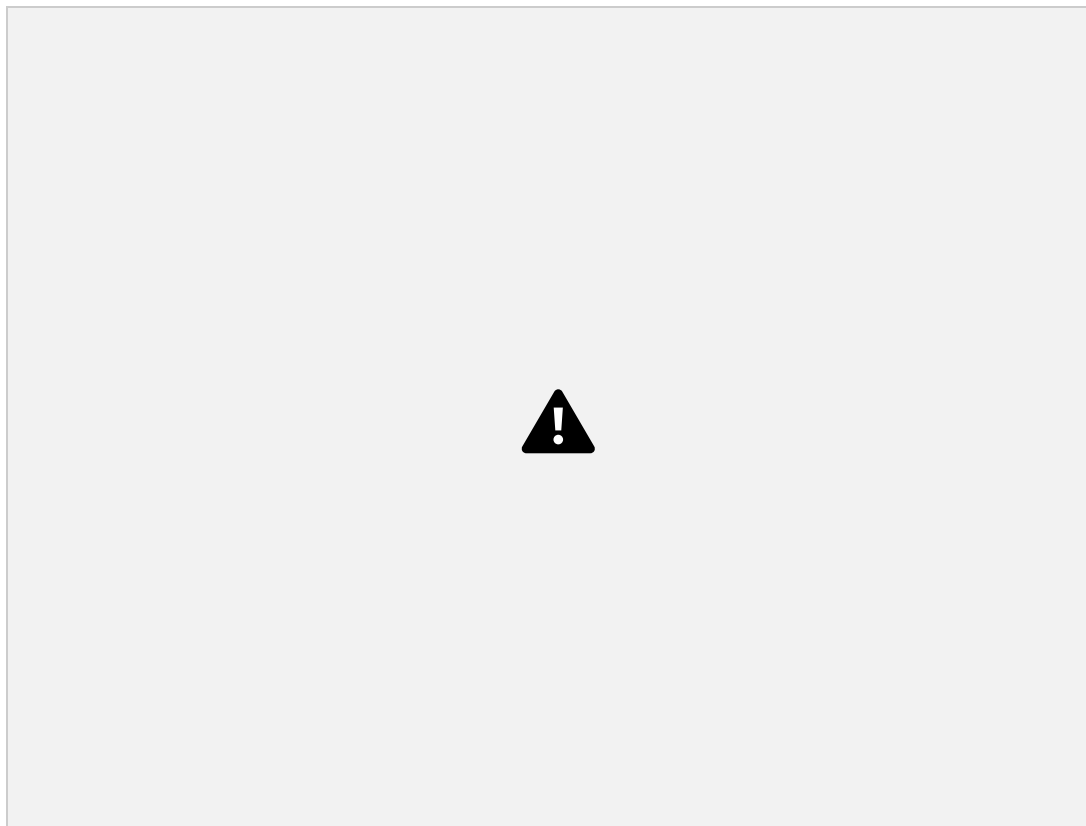Distribution of Waterpoint Type (False Positives vs. True Functional)

Identifying

top false positive features

Feature engineering on top offenders

Accuracy score ultimately ends up around 81%



**Help bring clean water to Tanzania**

# Thank you!