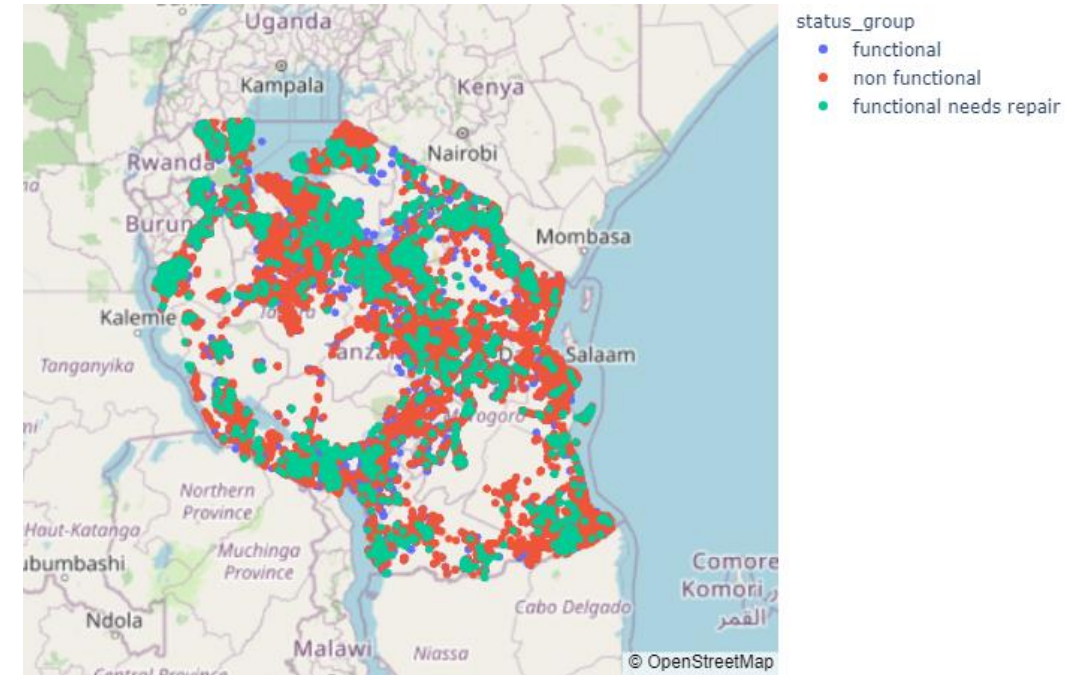# Classifying Functional and non-Functional wells

Rafael Villanueva

# Problem Statement



- Tanzanian's access to water is tied to accessibility to a working pump

- The Tanzanian Ministry of Water conducted a survey on several aspects of water pumps resulting in one outstanding question: are they still functional?

# Data Cleaning

# Modeling

# Findings
## Ensemble Classifier's Excel

| Method | Train Accuracy Baseline | Tuned Model Train Accuracy | Tuned Model Test Accuracy |
|---|---|---|---|
| Decision Tree | 72% | 74.20% | 76.86% |
| Random Forest | 78.92% | 79.91% | 80% |
| XGBoost Classifier | 84.69% | 80% | 80% |

# Model Selection
## XGBoost: With Reservations



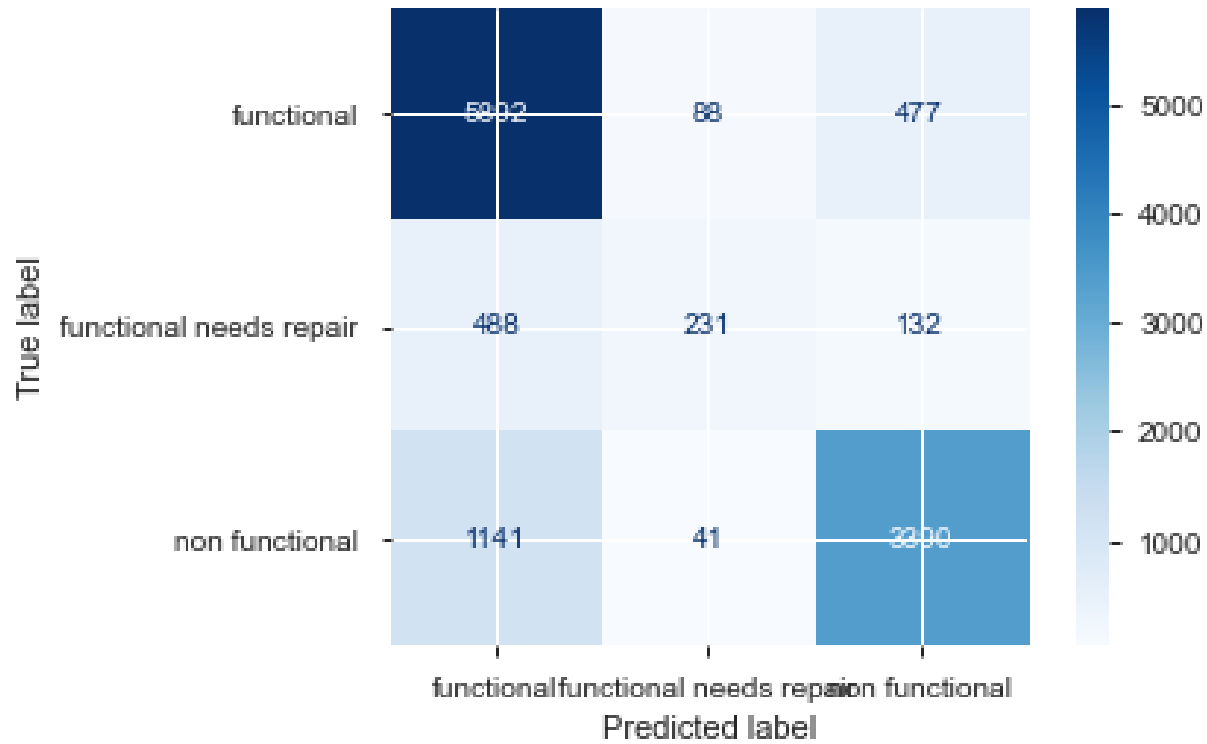|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| functional | 0.78 | 0.92 | 0.84 | 6457 |
| functional needs repair | 0.65 | 0.26 | 0.37 | 851 |
| non functional | 0.85 | 0.73 | 0.79 | 4572 |
| accuracy |  |  | 0.80 | 11880 |
| macro avg | 0.76 | 0.64 | 0.67 | 11880 |
| weighted avg | 0.80 | 0.80 | 0.79 | 11880 |

# Recommendations

- Complex models may not always be superior
- Predicting water pump functionality is possible!

# Future Work

- Focus on preprocessing data
  - In depth research on how the model operates on types of data, missing data etc.

# Thank You!

Eli Thomas

Flatiron Cohort 3.02.2020