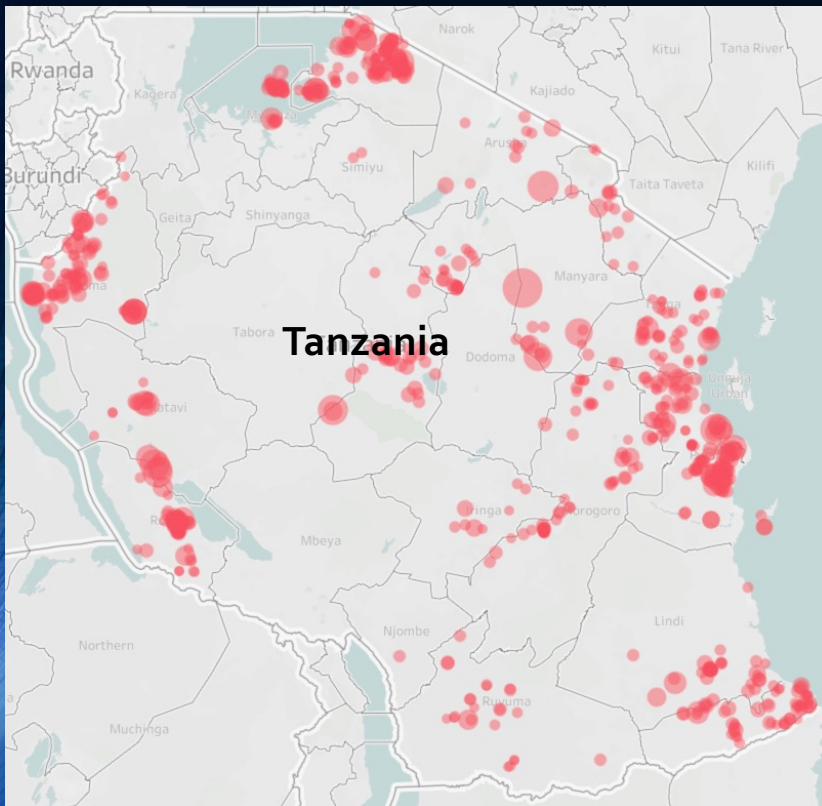# Predicting Water Well Failures Using Machine Learning

ANALYSIS BY BRIAN BENTSON

# Predicting Well Failure Can Save Lives

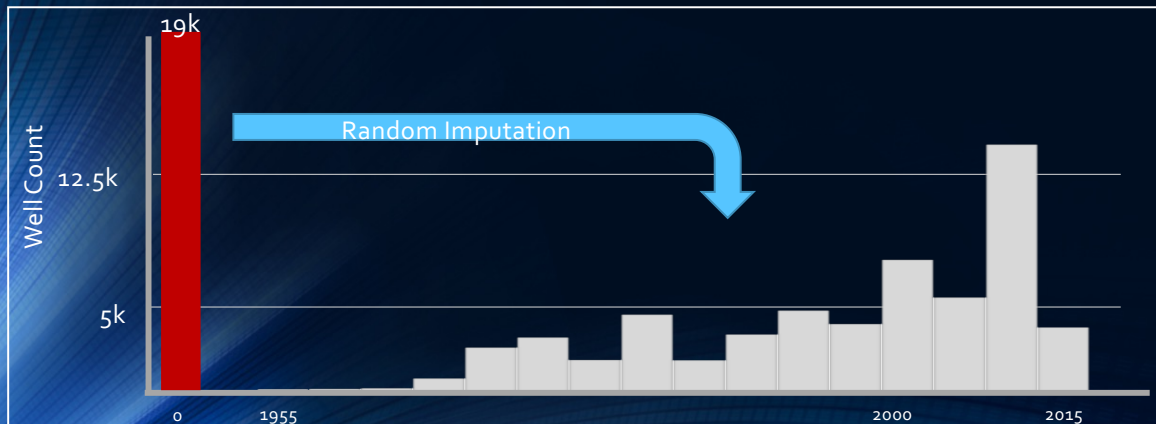Non-Functional wells that support over 1,000 citizens
Sized by population



- Tanzania has almost 60,000 water wells in dataset with 45% not functional, leaving almost 7M people without a reliable water source

- Humans can only live up to 3 days without water

- Ability to predict water well failures and respond quickly can be the difference between life and death

# Data Quality a Potential Issue

# 18K

- Almost 18,000 water wells with zeros for population, head, well elevation and construction year

- Recommendation: Improving data quality can drastically improve modeling performance

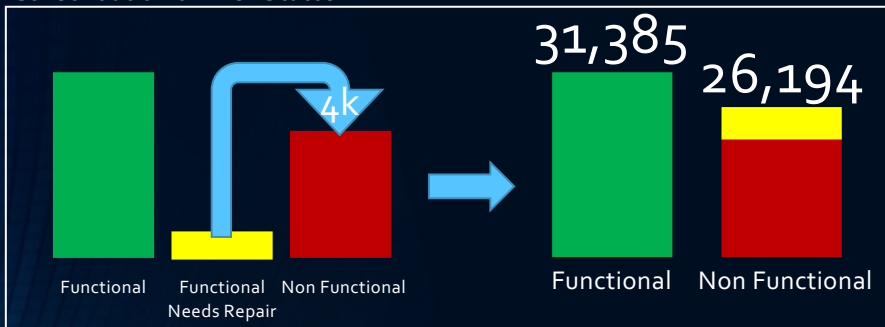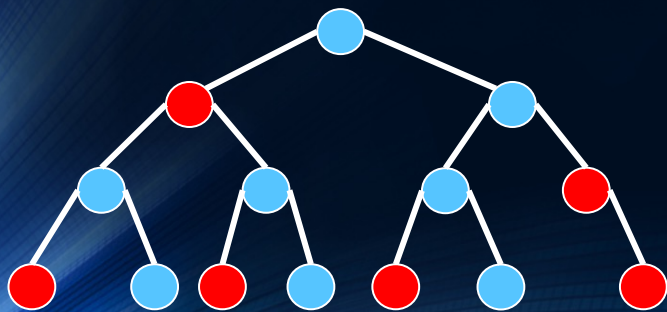## Distribution of Well Construction Year



- Utilized random imputation to convert zeros to non-zero values

- Outliers and improbable values in the head that can skew results

# Analysis Overview

### Consolidation of Well Status



## Random Forrest
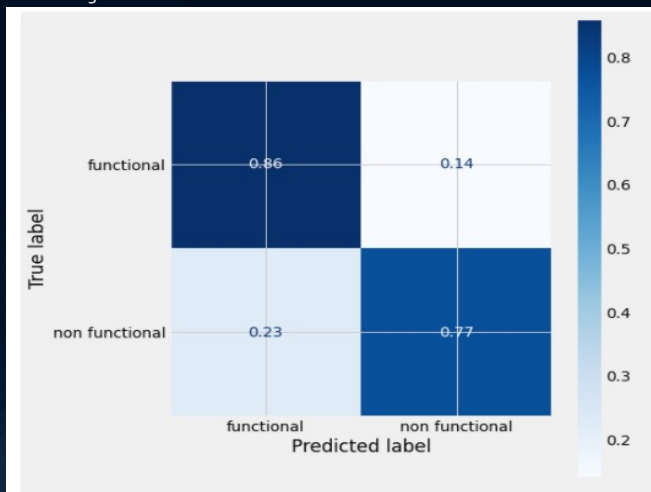


- Equated wells functioning at a reduced capacity as not functioning

- Prioritized finding all well failures, therefore accepting some false positives

- Focused on models with high interpretability to understand what drives well failures
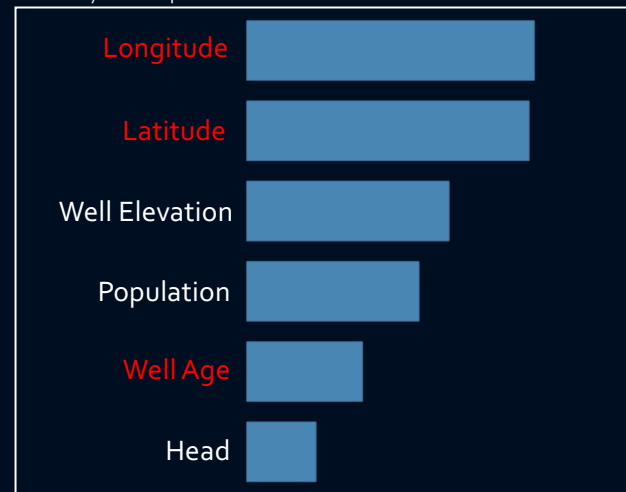
# Best Classification Modeling of Well Status

**Random Forest Performance**
Describing the final model



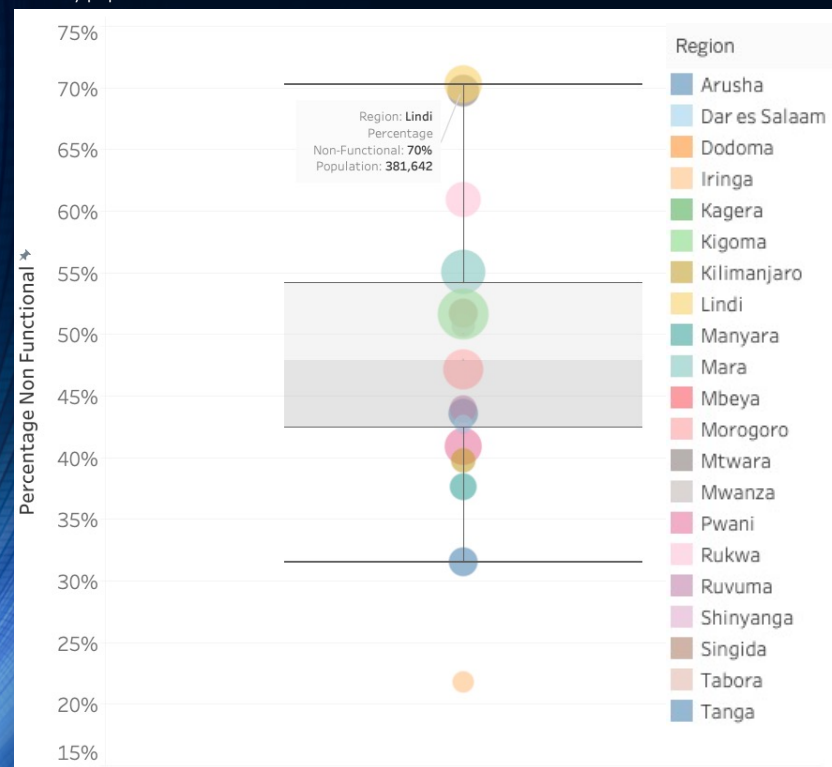**Random Forest Feature Importances**
Sorted by Most Important



- Best model was a Random Forrest that could detect failures 77% of the time with an overall accuracy of 82%

- Location, well age and head are important for the model

# Location Affects Reliability
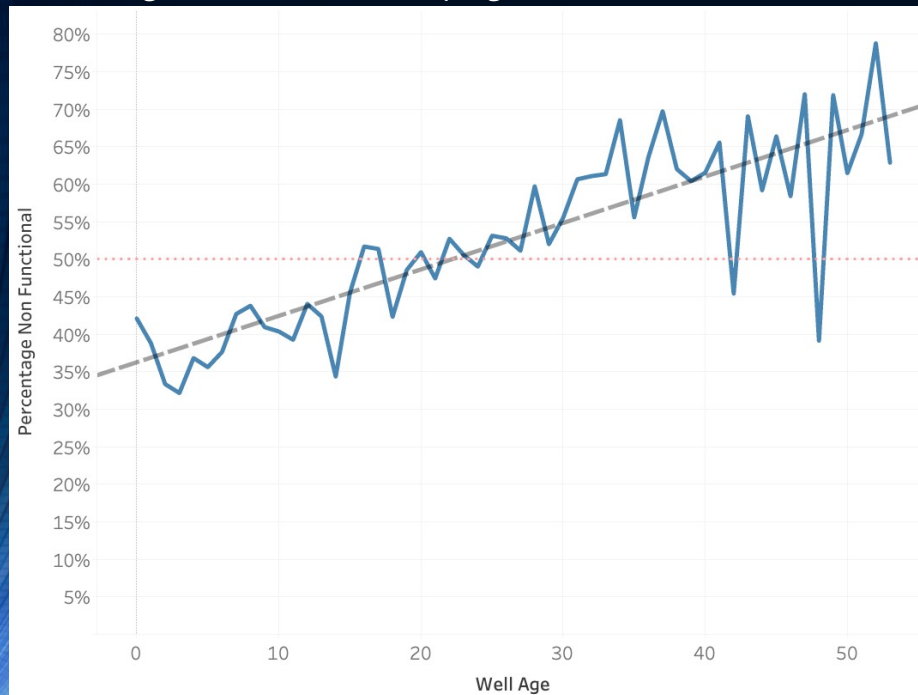
Distribution Percentage Non-Functional by Region
Sized by population



Region: **Lindi**
Percentage
Non-Functional: **70%**
Population: **381,642**

Region
- Arusha
- Dar es Salaam
- Dodoma
- Iringa
- Kagera
- Kigoma
- Kilimanjaro
- Lindi
- Manyara
- Mara
- Mbeya
- Morogoro
- Mtwara
- Mwanza
- Pwani
- Rukwa
- Ruvuma
- Shinyanga
- Singida
- Tabora
- Tanga

- High variability in well function across regions

-  Many Regions with high failure percentage and high population

- Recommendation: Focus on regions which have historically high failure percentage
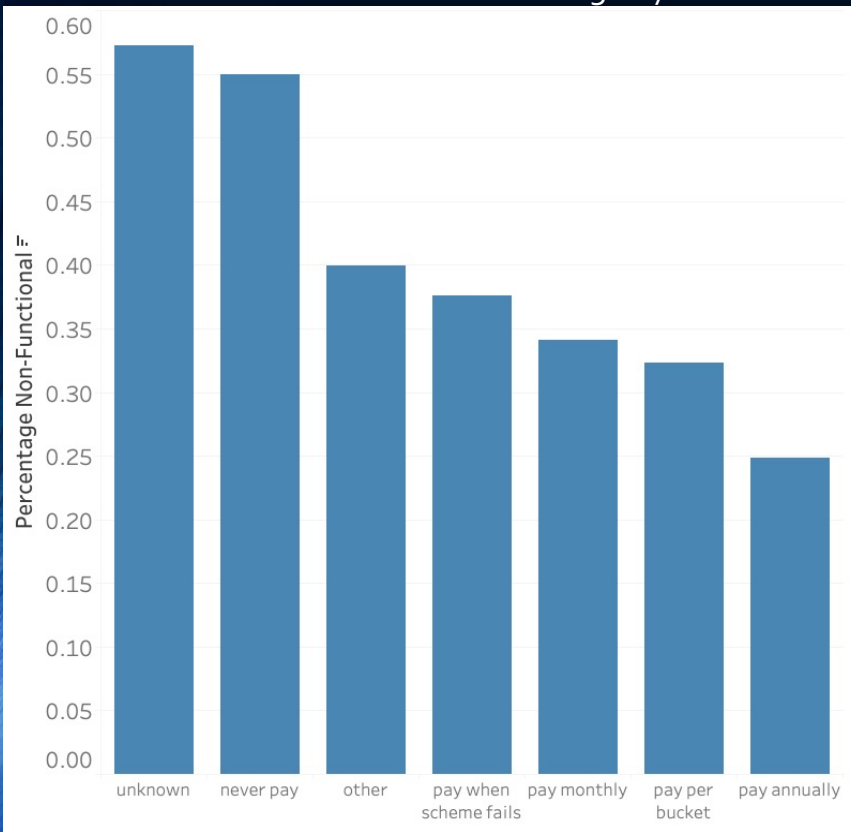
# Well Age Negatively Affects Reliability

Percentage Non-Functional by Age



- As well ages increase, the percentage of non-functional wells increases

- If a well is over age 24-25, it is more likely that well will be non-functional rather than functional

- 2.5M people are supported by older wells (>24 years old)

- Recommendation: Focus maintenance on older wells to maintain supply of water

# You Get What You Pay For

Distribution of Non-Functional Percentage by Water Cost



- While not specifically important for the random forest model, there is a clear trend between showing that if you pay for the water, the wells reliability is higher

- Recommendation: Focus on supporting the populations which cannot afford to pay for water
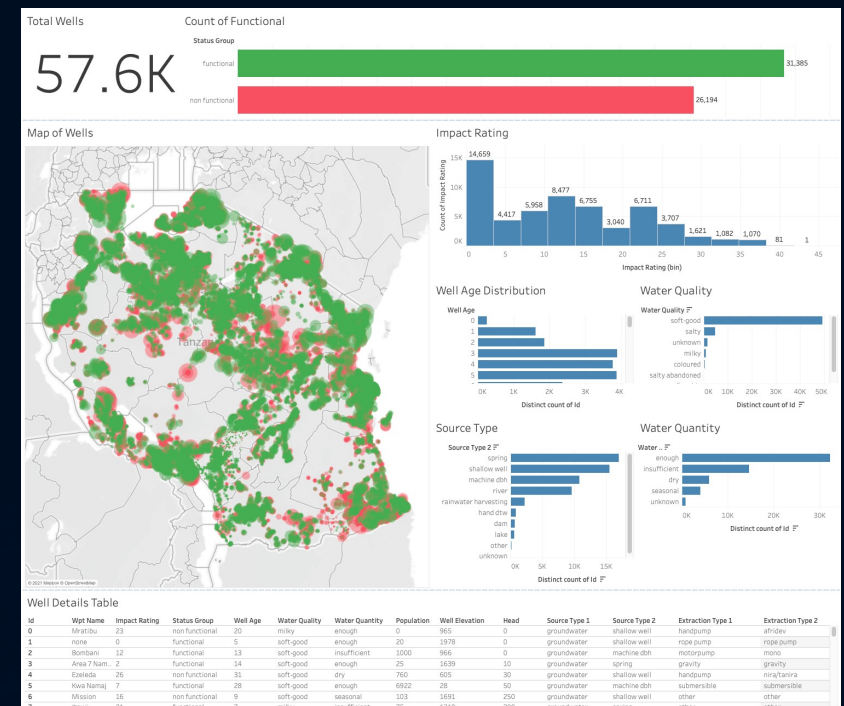
# Recommendations & Next Steps

## Recommendations

- Keep an eye on older wells as the likelihood of failure increases with age

- Develop support model for population areas that do not pay for water since analysis shows paying for water brings better reliability

- Utilize business insights tools to keep stakeholders up to date with key performance metrics

- Improve data governance to ensure better data quality and better predictions

## Next Steps

- Use more sophisticated algorithms that may perform better at finding well failures but does not tell you why they are failing
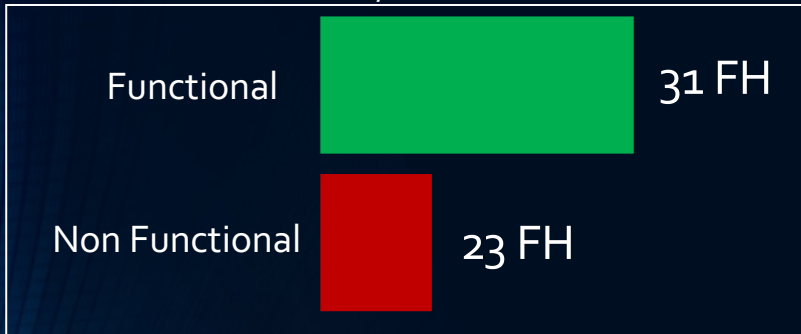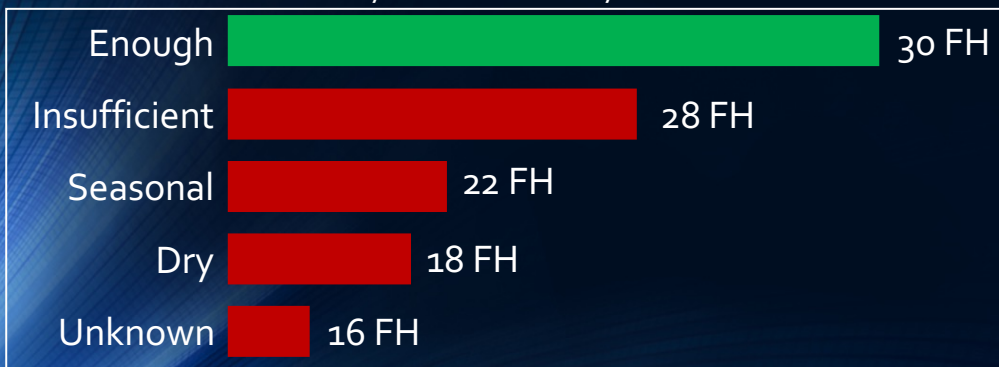
Screenshot of BI Dashboard

# Thank you

Questions?

# Appendix

# Lower Static Head Indicates Failure

Amount of Static Head by Well Status

| | |
|---|---|
| Functional | 31 FH |
| Non Functional | 23 FH |

Amount of Static Head by Water Quantity

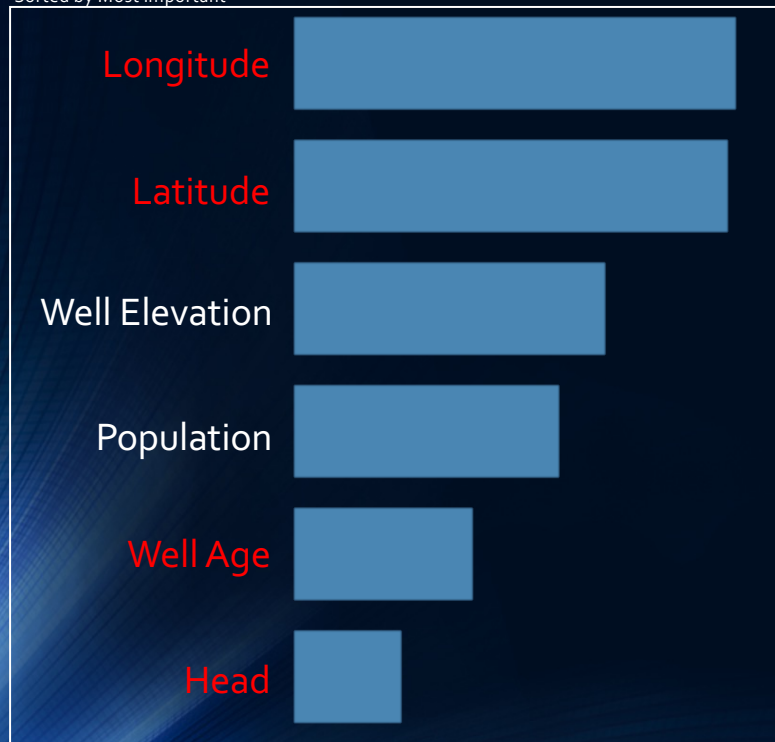| | |
|---|---|
| Enough | 30 FH |
| Insufficient | 28 FH |
| Seasonal | 22 FH |
| Dry | 18 FH |
| Unknown | 16 FH |

- As well ages increase, the amount of static head on the well will decrease, lowering water quantity

- This can be artificially improved by technology such as a pump

- Recommendation: Keep a close eye on static head as it directly correlates with water quantity

# Most Important Features from Model
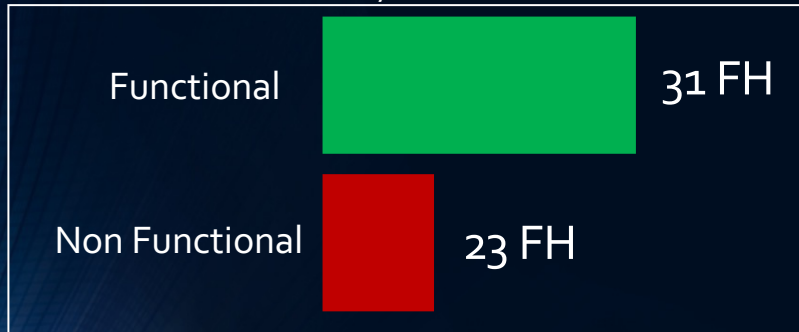
Random Forest Feature Importances

Sorted by Most Important



- Location (longitude and latitude) is the most important feature for predicting well function

- Well age directly affects reliability

- Well Elevation and Population are important features for prediction, although analysis did not highlight specific relationships
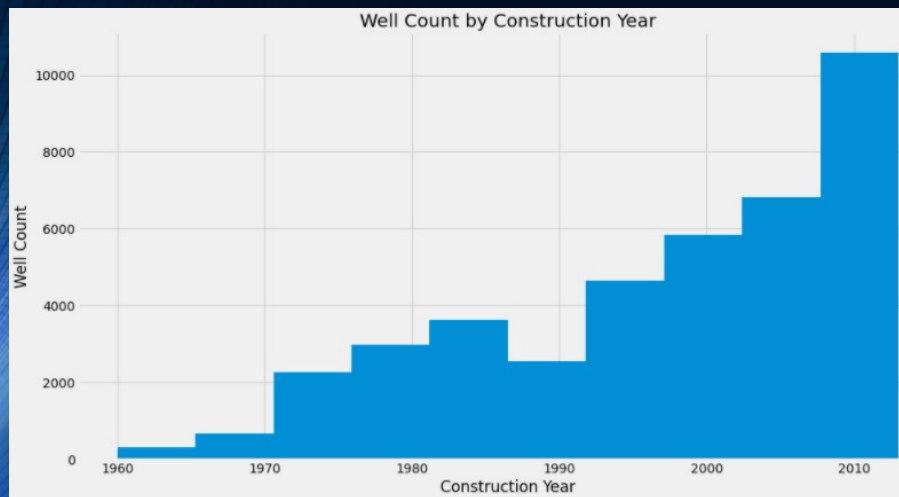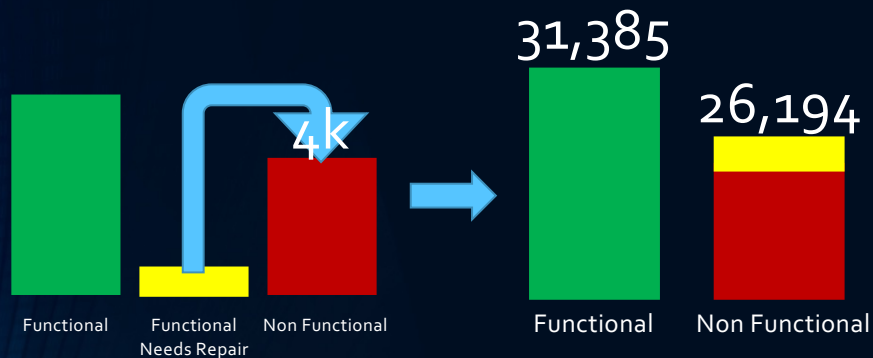
# Lower Static Head Indicates Failure

Amount of Static Head by Well Status

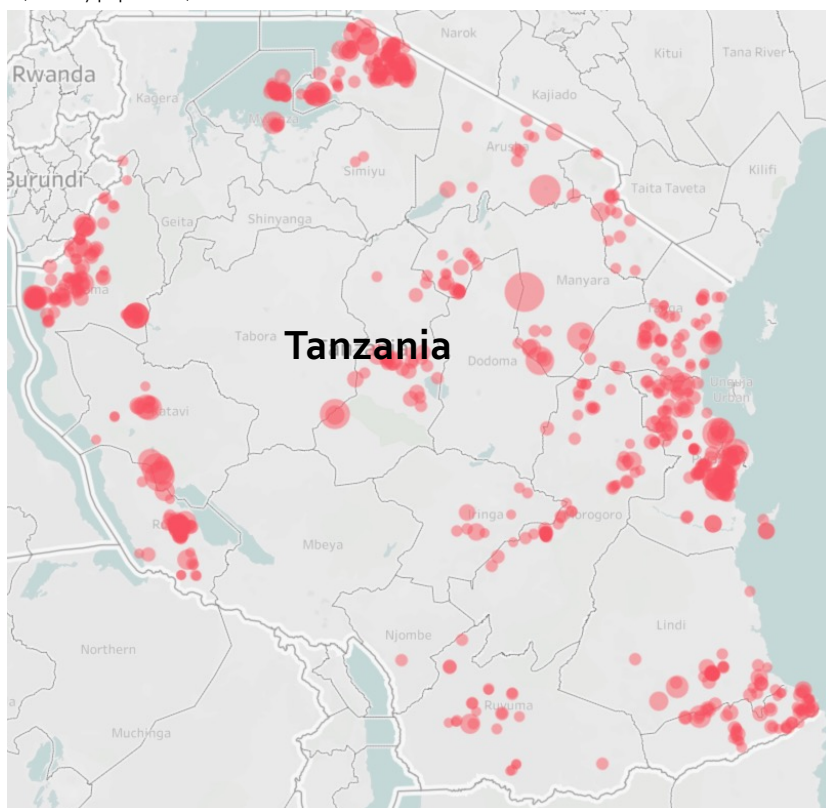Functional     31 FH

Non Functional     23 FH

- As well ages increase, the amount of static head on the well will decrease, lowering water quantity

- This can be artificially improved by technology such as a pump

- Recommendation: Keep a close eye on static head as it directly correlates with well function

# Analysis Overview



31,385

4k

26,194

Functional   Functional Needs Repair   Non Functional

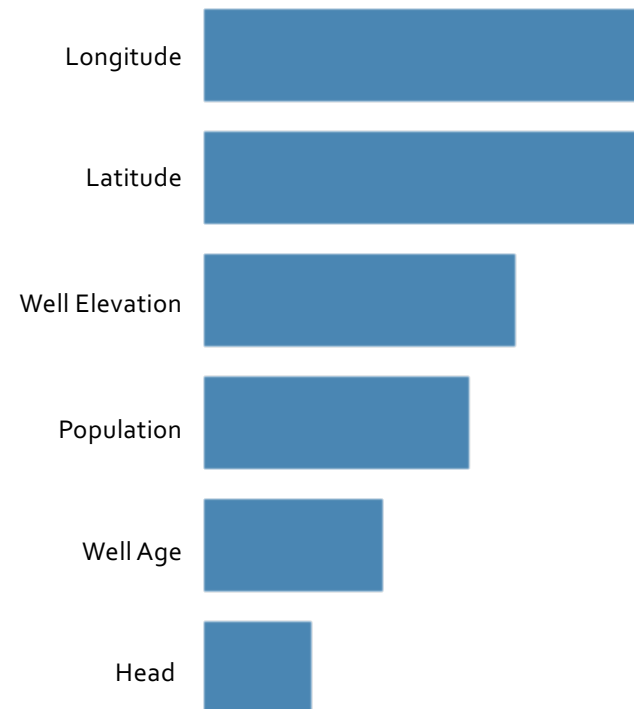Functional   Non Functional



Well Count by Construction Year

- Focused on wells not functioning or functioning at a reduced capacity (functional needs repair)

- Created well_age feature

- 20,000 wells with unknown construction_year. Filled values keeping identical distribution

- Classification modeling to use well features to predict non-functional water wells can save lives by increasing reliability and maintenance response time
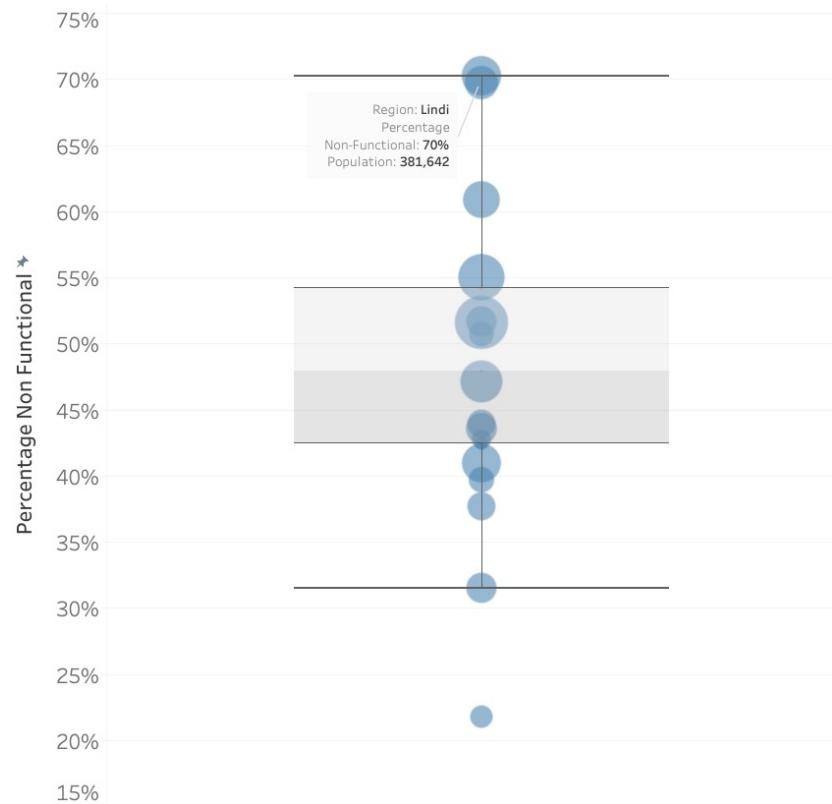
# Non-Functional wells that support over 1,000 citizens

(sized by population)
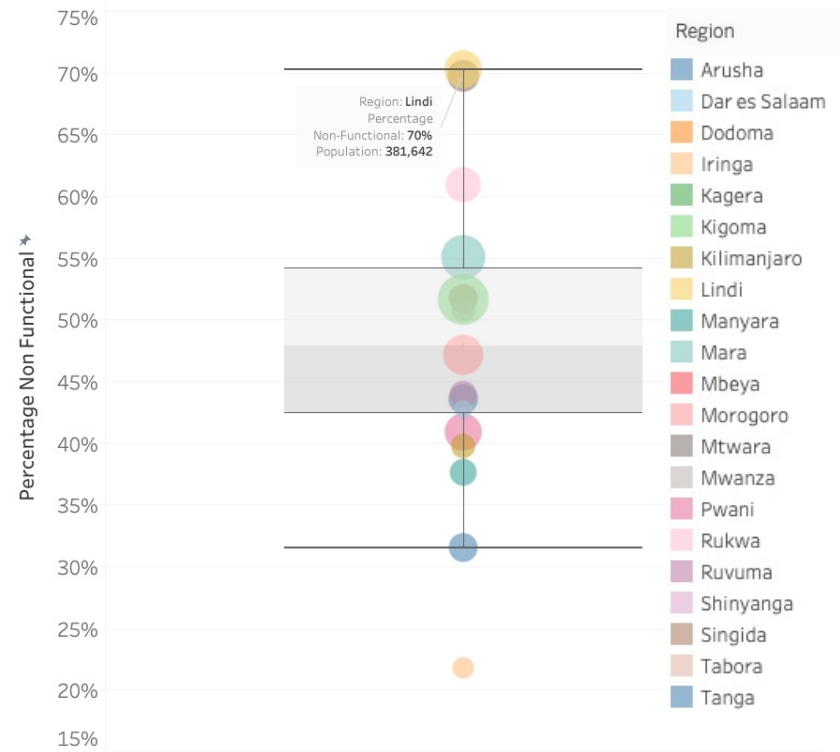
# Most Important Model Features

## Percentage of Non-Functional Wells by Region



Region: **Lindi**
Percentage
Non-Functional: **70%**
Population: **381,642**

Percentage Non Functional ✶

## Distribution Percentage Non-Functional by Region
Sized by population



Region: **Lindi**
Percentage
Non-Functional: **70%**
Population: **381,642**

Percentage Non Functional ✶

Region

- Arusha
- Dar es Salaam
- Dodoma
- Iringa
- Kagera
- Kigoma
- Kilimanjaro
- Lindi
- Manyara
- Mara
- Mbeya
- Morogoro
- Mtwara
- Mwanza
- Pwani
- Rukwa
- Ruvuma
- Shinyanga
- Singida
- Tabora
- Tanga

## Amount of Static Head by Well Status

Functional — 31 FH

Non Functional — 23 FH

## Amount of Static Head by Water Quantity

Enough — 30 FH

Insufficient — 28 FH

Dry — 22 FH

Seasonal — 18 FH

Unknown — 16 FH

# Percent Non-Functional by Well Age

# Random Forest Model Results

```
Classification Reports----------------------------------------
              precision    recall  f1-score   support

           0       0.82      0.86      0.84      9408
           1       0.82      0.77      0.80      7866

    accuracy                           0.82     17274
   macro avg       0.82      0.82      0.82     17274
weighted avg       0.82      0.82      0.82     17274

Test Graphs---------------------------------------------------
```