



Predicting Water Pump Operability in Tanzania

Claude Fried

October 2020



Problem Statement

How can we use information about these pumps to predict if they are working or not?

- The goal is to use data on water pumps in Tanzania to predict the condition of the well.
- There are thousands of water pumps in the country of Tanzania, many of which are an integral part of their area's sustainability.
- Many of these pumps are known to need maintenance (or are not working entirely).



Goal / Value

Goal:

- Our goal is to build a model that can accurately predict the condition of a water pump in Tanzania.

Value:

- Make data-driven decisions on how **to allocate resources** in order to maintain / fix pumps that are most in need.
 - **Save money** – spend only on pumps which are most in need of repairs.
 - **Save time** – avoid erroneously targeting pumps which are functional.
- **Improve maintenance operations** for these water pumps and significantly improve many people's quality of life.



Methodology – OSEMN

Obtain – Scrub – Explore – Model – iNterpret

1. Obtain: Gather data.
2. Scrub: Clean data.
3. Explore: Get to know the data. Engineer features.
4. Model: Train, test, validate, and compare models.
5. Interpret: Glean insights from the model's predictions.

Methodology – OSEMN

Obtain

- The data has been collected from Taarifa (<http://taarifa.org/>) and the Tanzanian Ministry of Water.
- 59,000 training data points.
- 41 original features.



Methodology – OSEMN

Scrub

Examples:

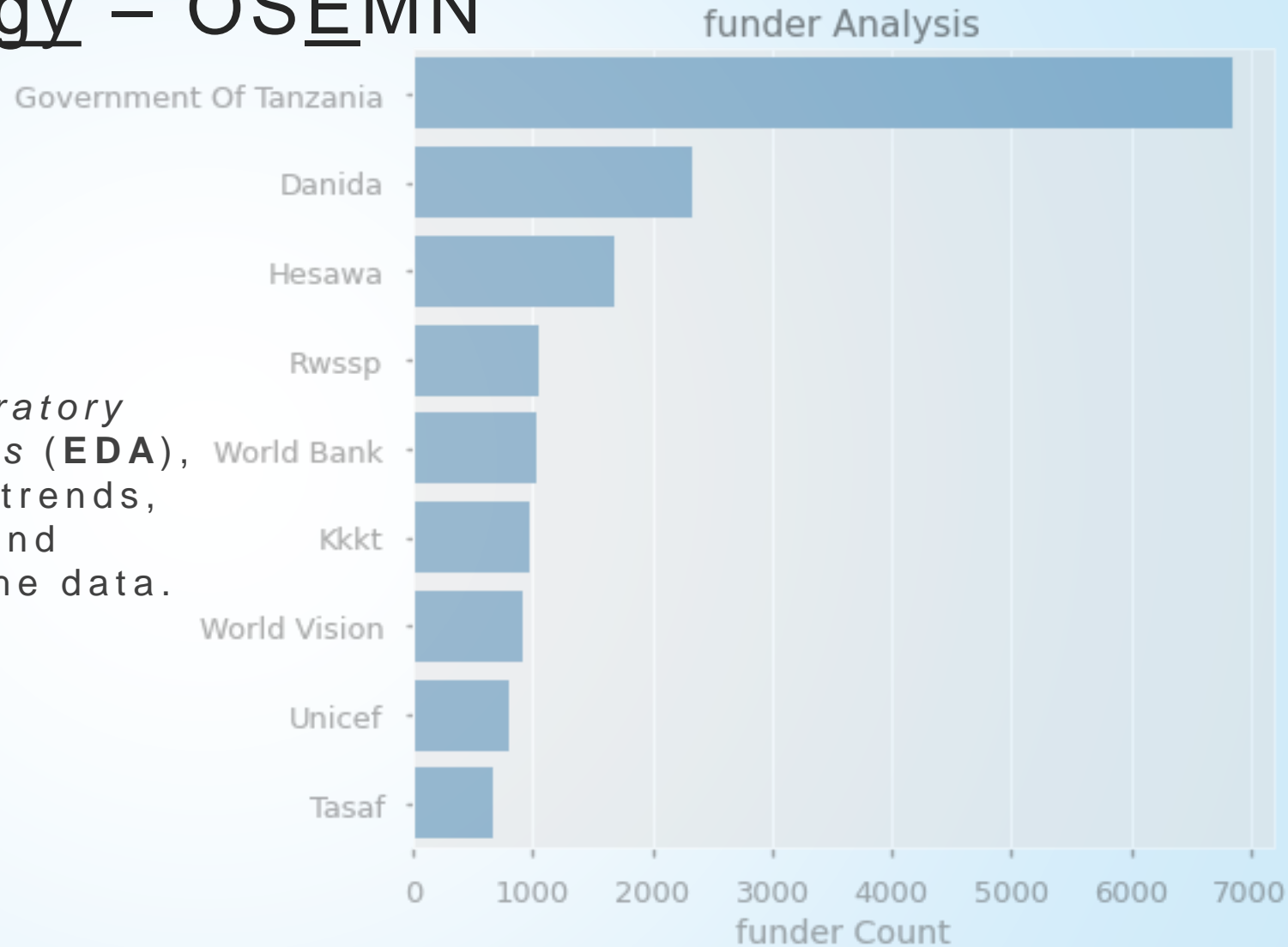
- Format each feature into the correct data-type.
 - e.g.: Is this feature numerical or categorical?
- Determine what to do with...
 1. Missing values for each feature.
 2. Numerous categories within a feature.
 3. Duplicate / similar features.



Methodology – OSEMN

Explore

- Using *exploratory data analysis (EDA)*, we can find trends, groupings, and outliers in the data.





Methodology – OSEMN

Model

The modeling process included several steps.

1. Selecting which model-types to train.
2. Searching for optimal settings for these models.
3. Creating a new model made up of these trained and optimized models.
4. *Select a final model from which we will gain insights.*

Methodology – OSEMN

Final Model – Random Forest

The model that best-balanced performance and speed was the **Random Forest Classifier**.

Random Forests use a collection of *Decision Trees* to make predictions.



Decision Trees

Tree_0

- extraction_type_group gravity <= 0.5
samples = 100.0%
value = [0.543, 0.073, 0.384]
 - waterpoint_type_group other <= 0.5
samples = 54.9%
value = [0.497, 0.049, 0.454]
 - longitude <= -2.926
samples = 46.1%
value = [0.571, 0.054, 0.375]
 - quantity_dr samples = 1.4%
value = [0.569, ...]
 - source_machine samples = 1.4%
value = [0.571, ...]
 - Iga_Siha <= 0.5
samples = 1.4%
value = [0.103, ...]
 - public_meeting_samples = 1.4%
value = [0.571, ...]
 - lon sa samples = 1.4%
value = [...]
 - scheme_management_W samples = 1.4%
value = [0.726, ...]
 - longitude <= 0.5
samples = 1.4%
value = [0.559, ...]
 - quantity_dry <= 0.5
samples = 1.4%
value = [0.127, 0.146, 0.728]
 - funder_Private <= 6.967
samples = 8.8%
value = [0.107, 0.026, 0.867]
 - gps_height <= 1.01
samples = 43.1%
value = [0.615, 0.099, 0.286]
 - gps_height <= -0.941
samples = 2.0%
value = [0.256, 0.139, 0.605]
 - waterpoint_type_group other <= 0.5
samples = 45.1%
value = [0.599, 0.101, 0.3]

Tree_99

- amount_tsh <= -0.131
samples = 100.0%
value = [0.543, 0.073, 0.384]
 - extraction_type_group gravity <= 0.5
samples = 70.2%
value = [0.475, 0.073, 0.452]
 - management_wug <= 0.5
samples = 40.7%
value = [0.436, 0.053, 0.511]
 - quantity_dry <= 0.5
samples = 29.5%
value = [0.529, 0.101, 0.371]
 - Iga_Kigoma Rural <= 4.174
samples = 3.2%
value = [0.474, 0.278, 0.248]
 - water_quality_salty <= 0.5
samples = 26.6%
value = [0.732, 0.047, 0.221]
 - longitude <= -0.402
samples = 29.8%
value = [0.704, 0.072, 0.224]
- waterpoint_type samples = 1.4%
value = [0.4, 0. ...]
- district_code samples = 1.4%
value = [0.585, value = [0.593, ...]
- longitude <= 0.5
samples = 1.4%
value = [0.593, ...]
- funder World Ba samples = 1.4%
value = [0.017, (value = [0.53, ...]
- funder Dwe samples = 1.4%
value = [0.53, ...]
- payment_type m samples = 1.4%
value = [0.333, value = [0.745, (value = [0.576, 0.06, 0.364]
- quantity_enou samples = 1.4%
value = [0.745, (value = [0.576, 0.06, 0.364]
- quantity_dry <= 0.5
samples = 2.0%
value = [0.576, 0.06, 0.364]

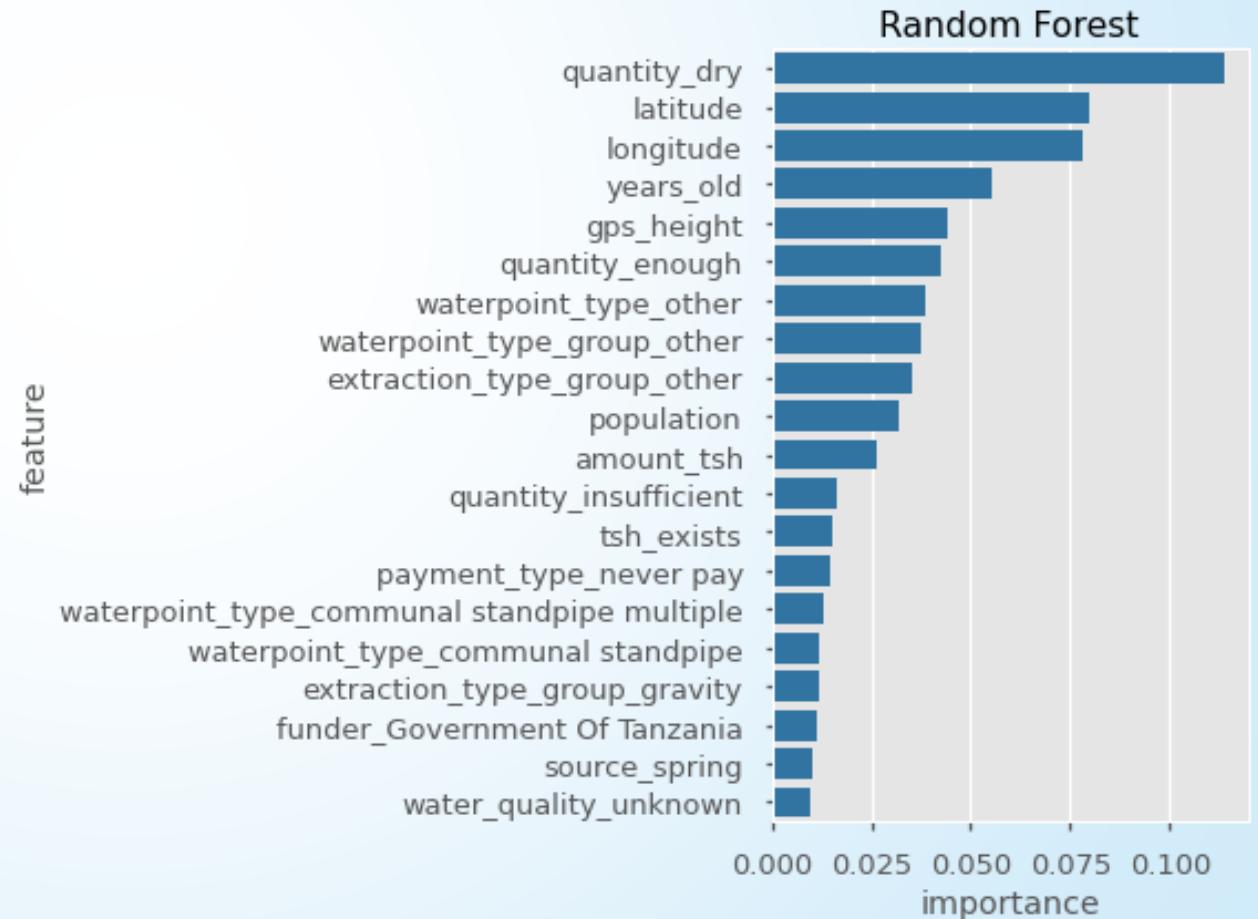
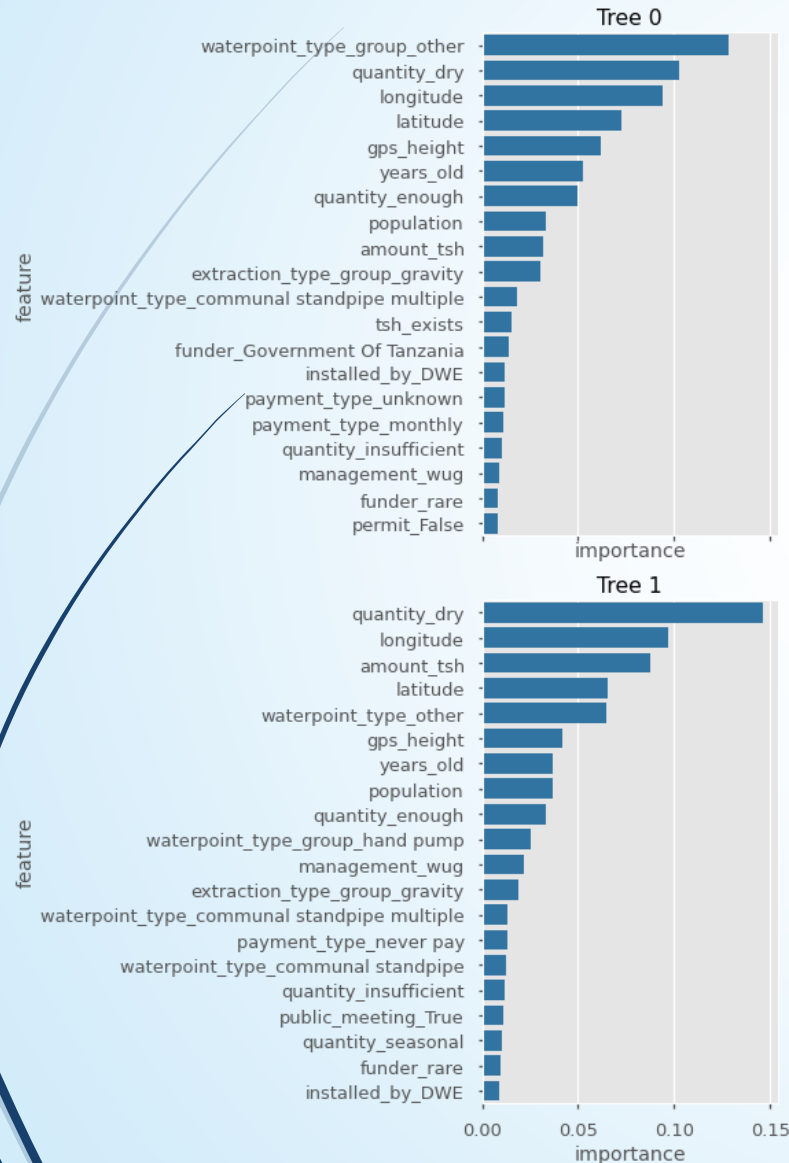
The Forest that was used to make final predictions and insights uses 100 unique trees which split the data (over 44,000 points) until each split contains a single entry.

That's a huge tree!

Methodology – OSEMN

Decision Trees

Decision Trees and Random Forests can show the most important features that they used to make their predictions.





Methodology – OSEMN

Interpret

- ▶ One of the great advantages of a Random Forest Classifier is its **interpretability**.
- ▶ Because of its method of classification (using splits of the data), the most vital features can be easily obtained.
 - ▶ A Random Forest using ***all 137 features*** of the data scores **81% accuracy**.
 - ▶ A Random Forest using ***only the 11 most-important features*** scores **79% accuracy**.

Methodology – OSEMN

Interpret

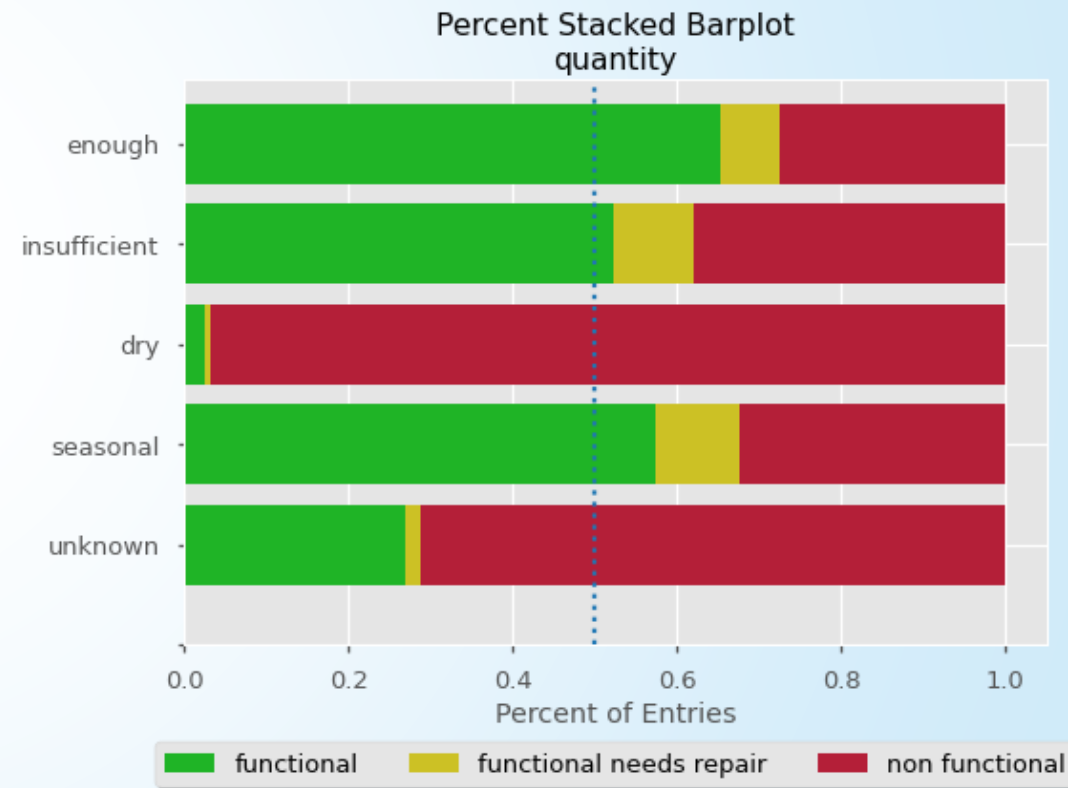
The 5 most important features for the Forest are:

1. Quantity - dry: Is the well dry or not?
2. Latitude: Geographical latitude of the well.
3. Longitude: Geographical longitude of the well.
4. Years old: How many years ago was the well installed?
5. GPS height: Altitude of the well.

How can I tell if a well is functional?

Recommendations

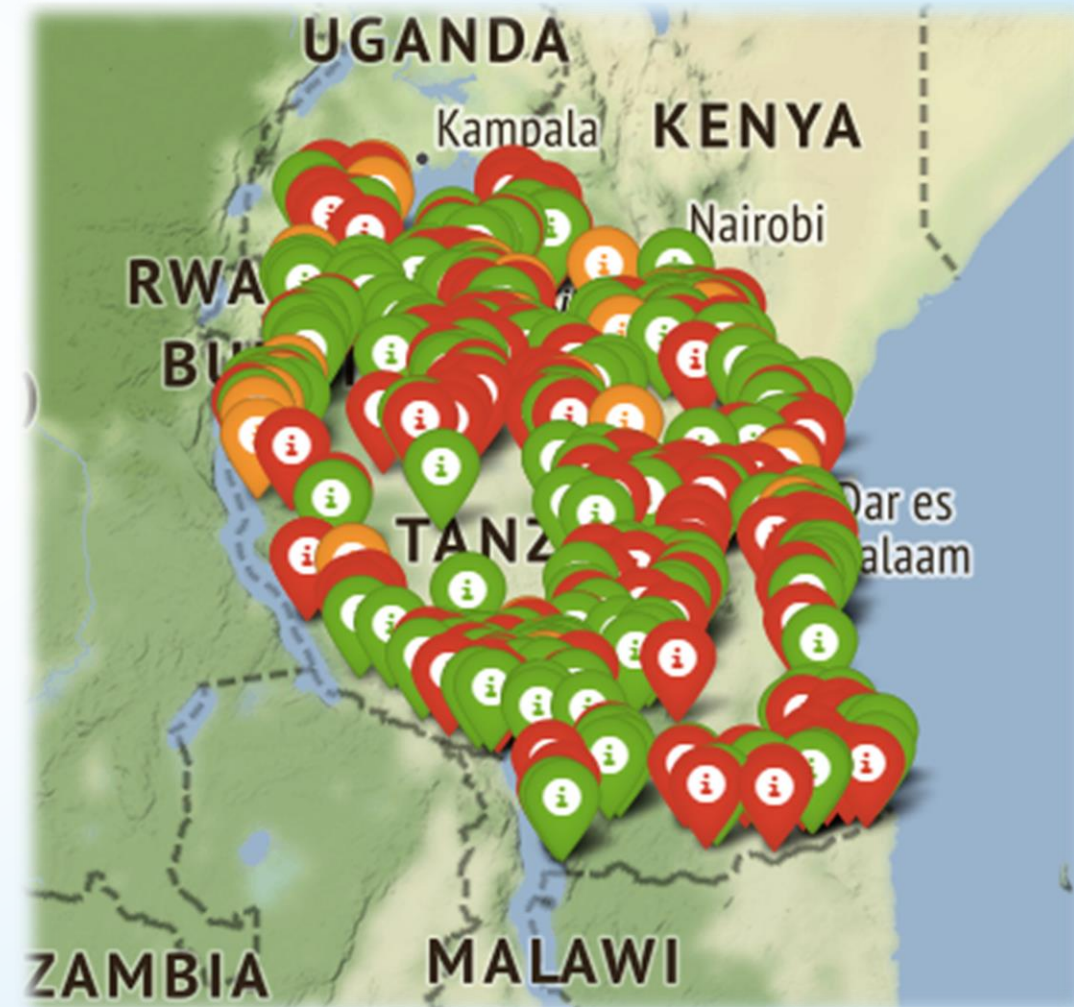
1. **Quantity:** Look for wells that with dry or unlabeled water quantities.
 - Wells that are labeled as **dry** are almost all *nonfunctional*.
 - Conversely, less than 25% of wells that have **enough** water are *nonfunctional*.
 - Note: The missing values for this feature are not missing at random (i.e.: when the label is **unknown**, it is much more likely to be *nonfunctional*).



How can I tell if a well is functional?

Recommendations

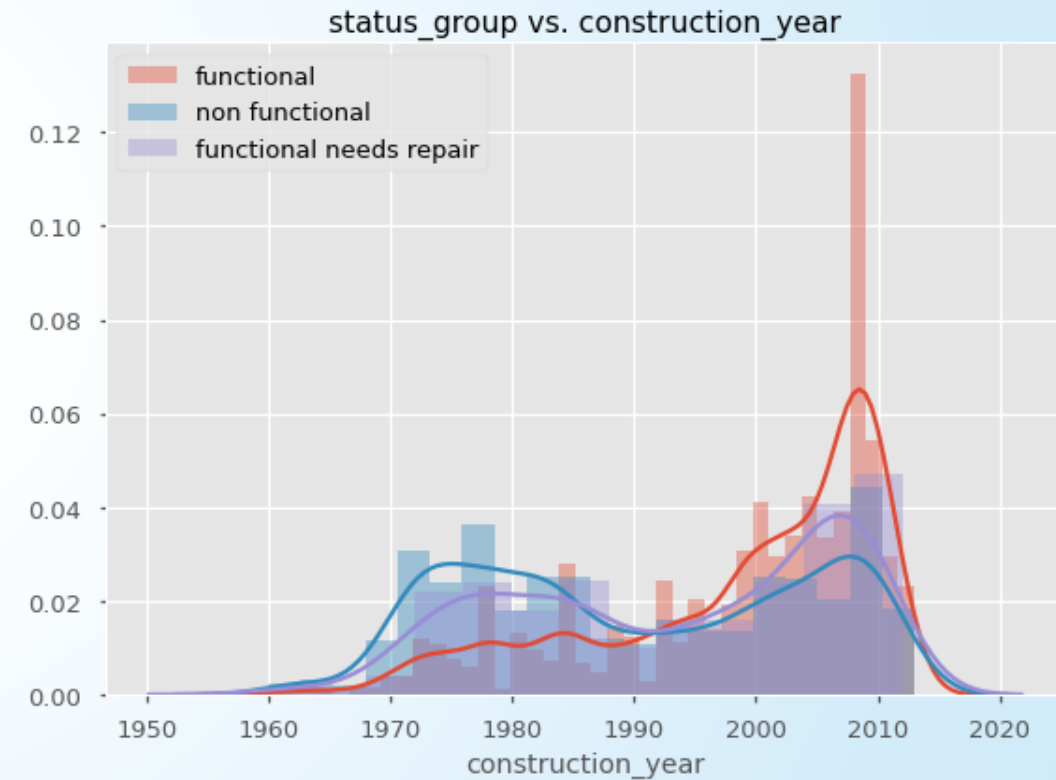
2. **Latitude / Longitude:** While there are functional and *nonfunctional* wells spread over the country, there do seem to be "pockets" of non-functional wells. This could be due to the water source or other geographical features of an area.
 - If there is a pocket of *nonfunctional* wells in an area, other wells in the same area might be *nonfunctional* as well.



How can I tell if a well is functional?

Recommendations

- 3. Years Old:** Look for older wells (especially built before 1990).
 - As the construction year increases, the wells become much more likely to be functional.
 - Wells built before **1990** are more likely to be *nonfunctional* or *needing repair*, while after 1990 are more likely to be *functional*.

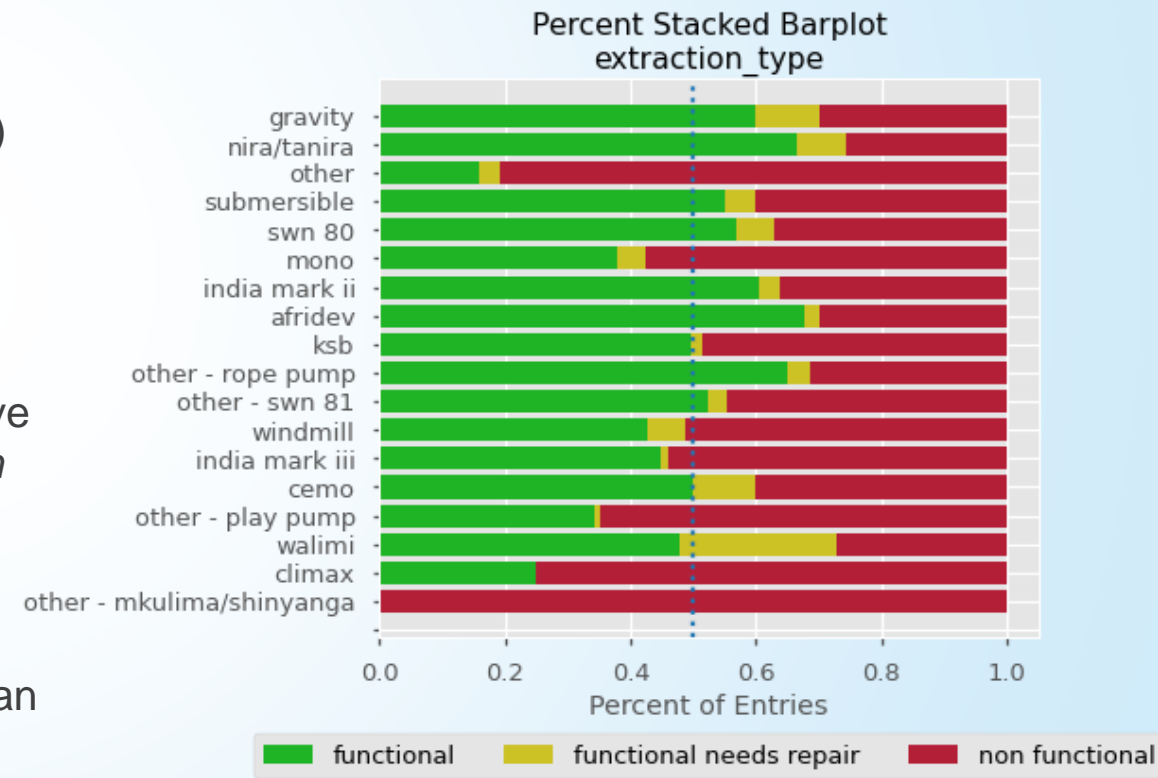


How can I tell if a well is functional?

Recommendations

4. Extraction Type: Determine the extraction type.

- ▶ **Gravity** pumps (the most common) are 60% likely to be *functional* and 25% to be *nonfunctional*.
- ▶ The most common type or well to have more *nonfunctional* wells than *functional* is **mono**. These have over a 60% chance to be at least *in need of repair*.
- ▶ Note: Wells that are labeled **other** have a much higher likelihood of being *nonfunctional* than otherwise.

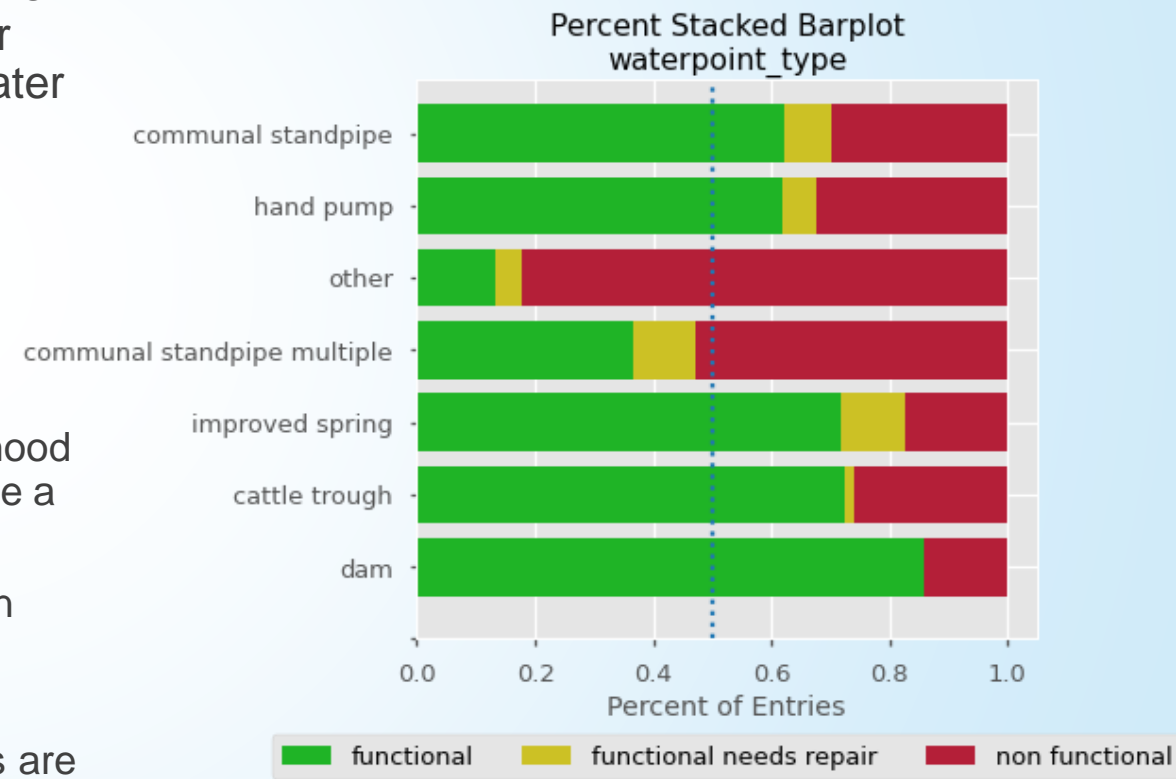


How can I tell if a well is functional?

Recommendations

5. **Waterpoint Type:** Find wells that have less common waterpoint types (other than a communal standpipe) for greater likelihood that it will need repairs.

- ▶ The most common **communal standpipe** is over 60% likely to be *functional* and 25% to be *nonfunctional*.
- ▶ Having a **communal standpipe multiple** though increases the likelihood of the pump being *nonfunctional* quite a lot.
- ▶ If the waterpoint type is less common (marked **other**), it is almost certainly (over 80%) *nonfunctional*.
- ▶ It seems that **improved spring** wells are the most likely to need repairs (these waterpoint types may require more frequent maintenance than others).



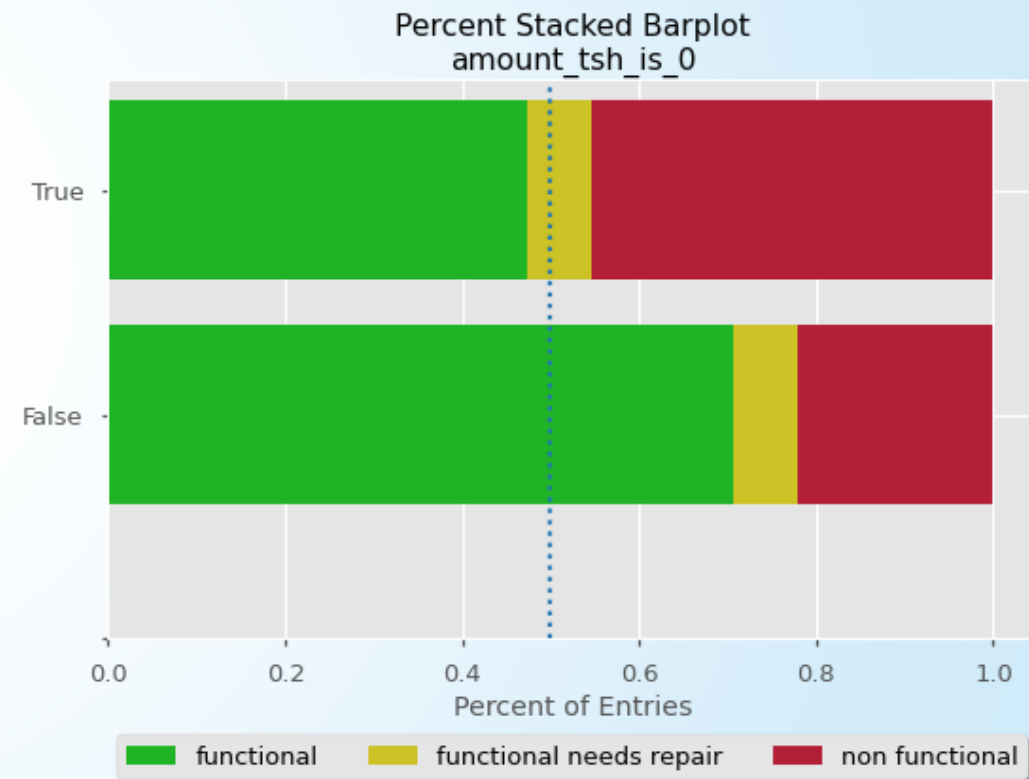
How can I tell if a well is functional?

Recommendations

- 6. Amount Total Static Head:** Most wells have 0 total static head (available static water), but some have as much as 200,000 units.

Wells that have 0 **tsh** are much more likely to be *nonfunctional* compared to wells with more than 0 **tsh**.

- ▶ As a general principal, the greater the **tsh**, the more likely it is to be *functional*.





Next Steps

- ▶ Dive deeper into the **location** of each well.
 - ▶ Are there regions where certain well-types work best?
 - ▶ Are certain water-sources likely to run dry soon?
 - ▶ Are there geographical features (mountains, deserts, plains) that impact the wells' condition?
 - ▶ Are wells within a close vicinity to one another more likely to be nonfunctional?



Thank You!

- ▶ Flatiron School
- ▶ Taarifa (<http://taarifa.org/>) and the Tanzanian Ministry of Water.