

# DataDrip

## Exploratory Data Analysis (Sprint 1)

Shakeel Ahmed  
Amey Chitnis

# Table of Contents

Introduction

Key Insights

- Dataset imbalance
- Construction Year
- Pumps per Region
- Initial columns exploration

# Introduction

Source: Taarifa & Tanzanian Ministry of Water

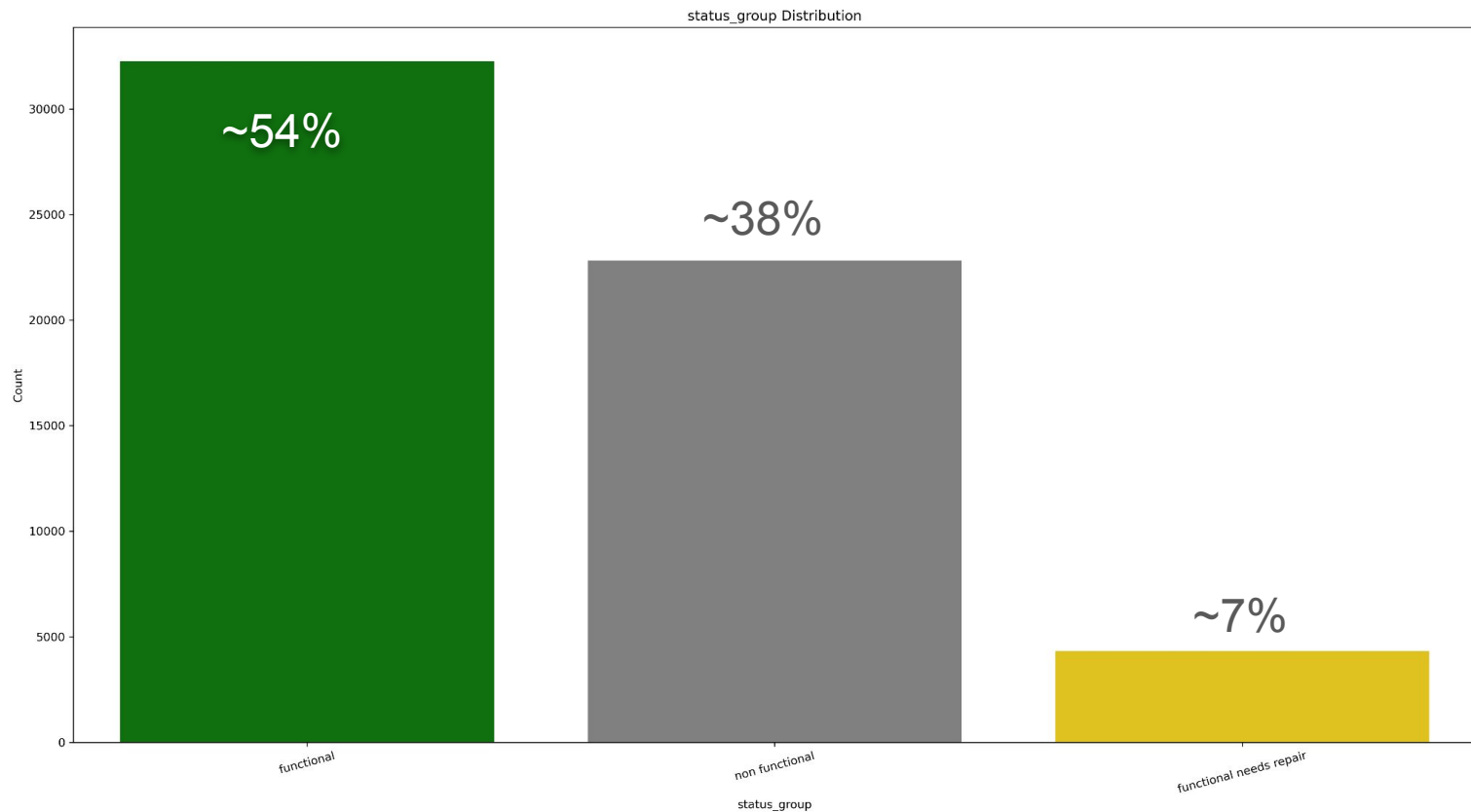
59000 water pumps

geolocation, management, water and temporal information

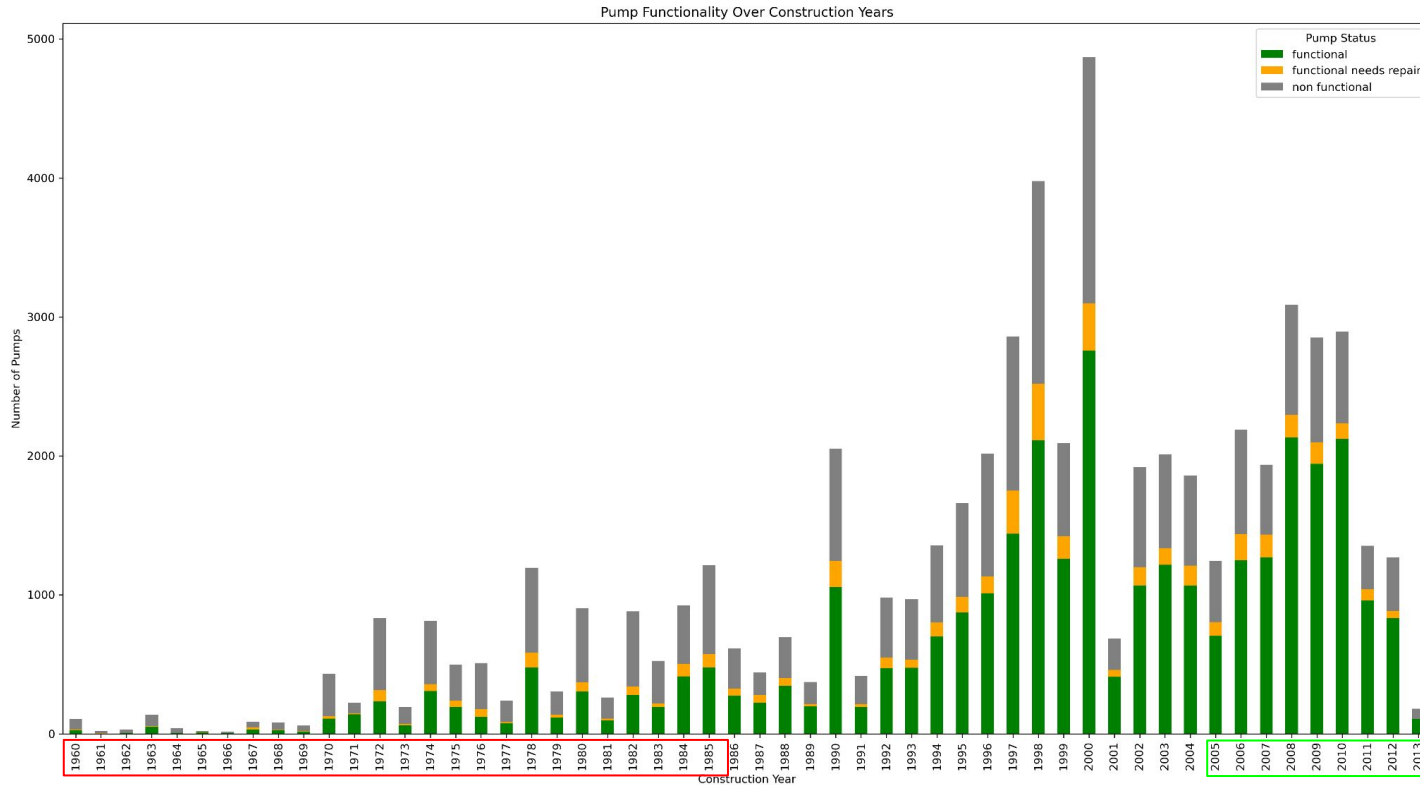
Functional status- Working, not working, needs repair

Several missing values: NaN, 0s and 1s

# Dataset Imbalance!

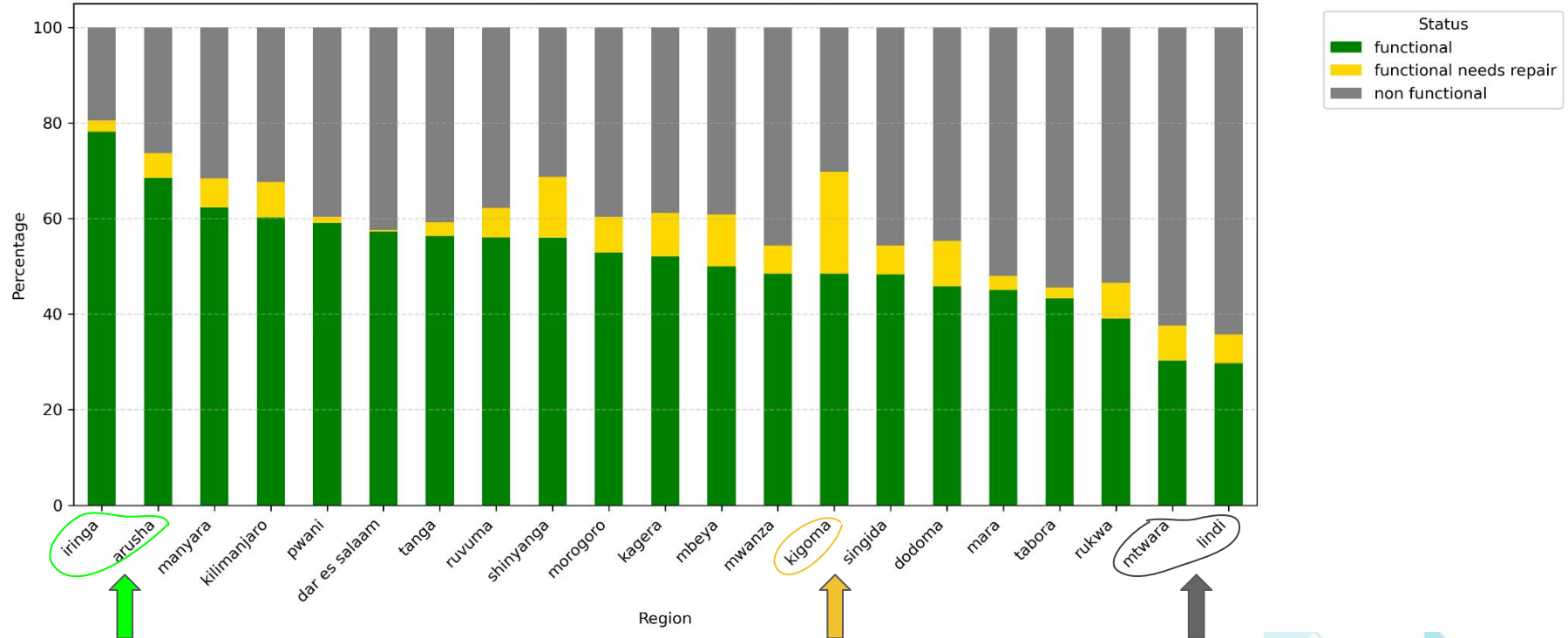


# Construction year

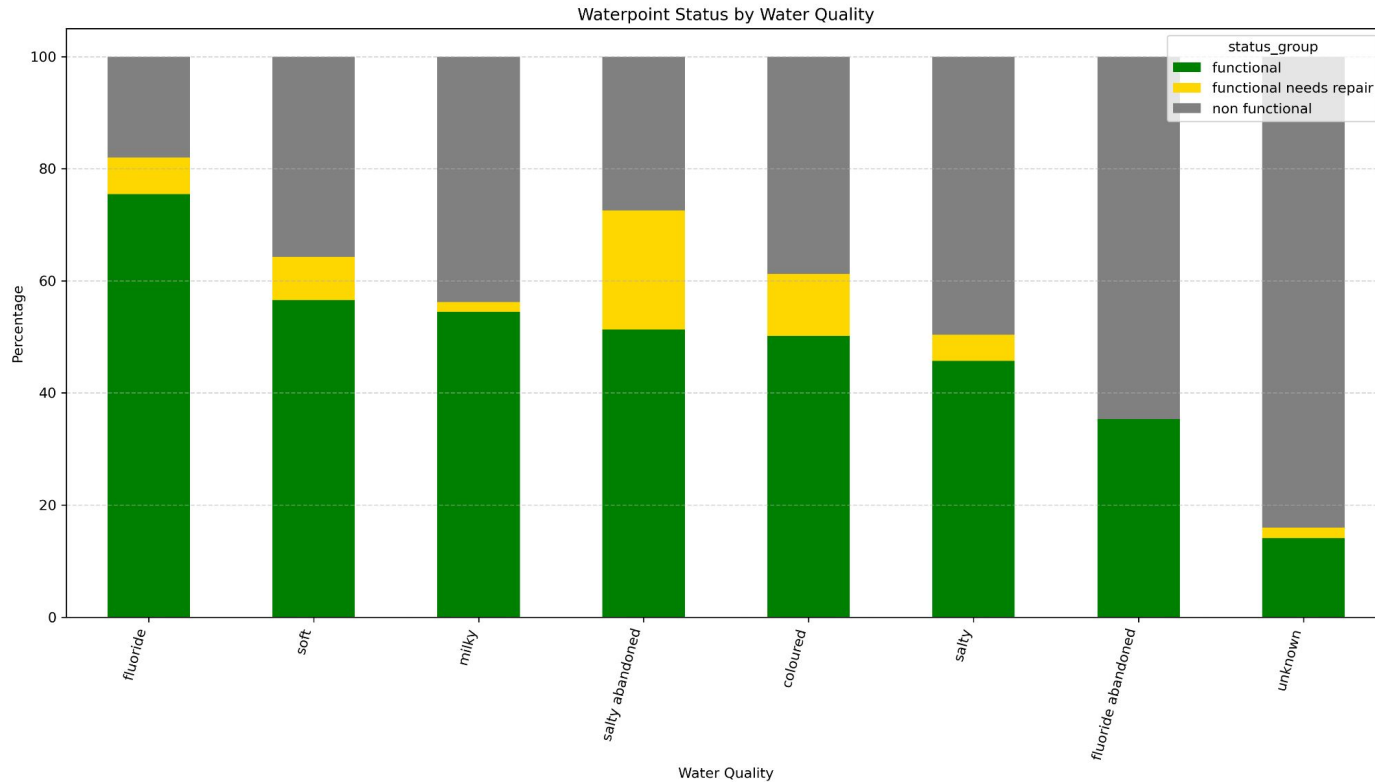


# Functional pumps per region

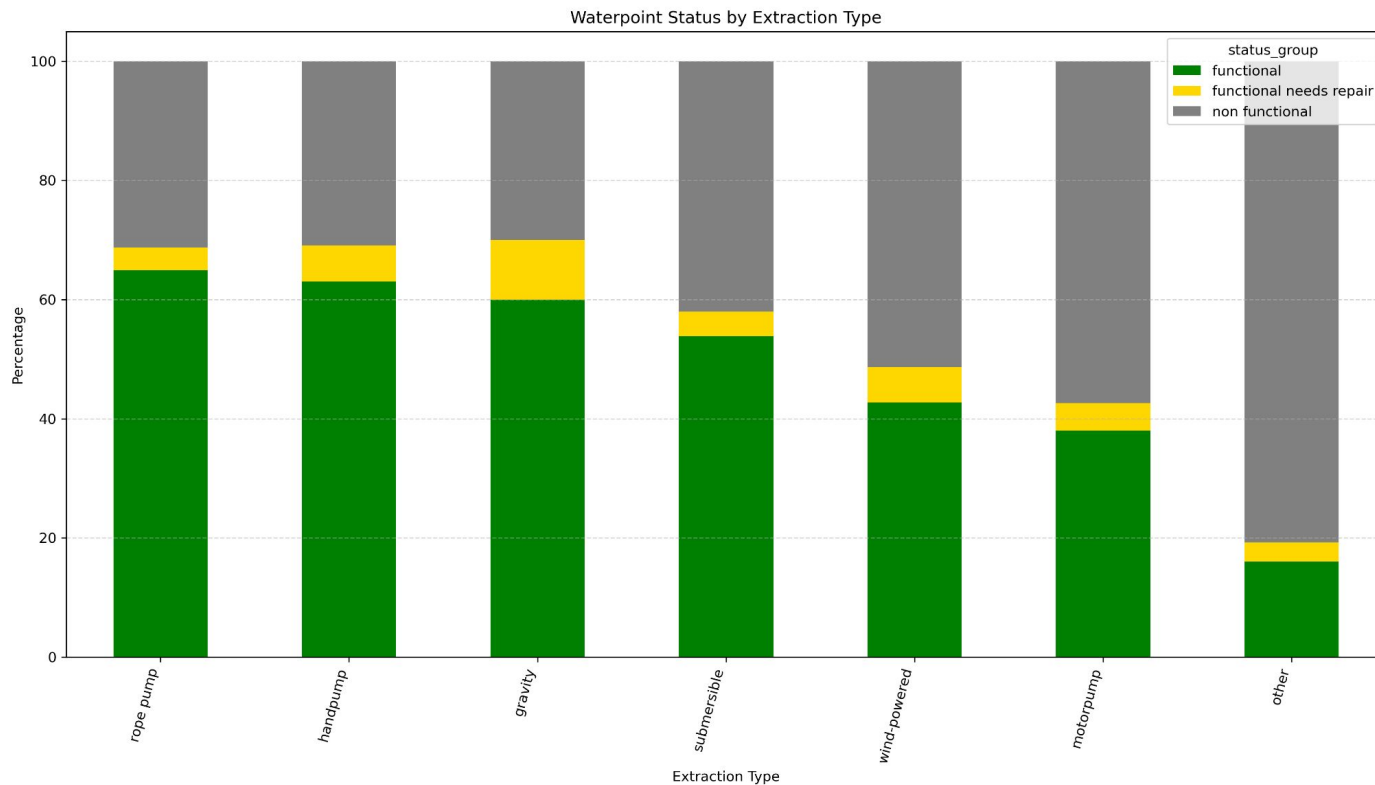
Waterpoint Status Distribution by Region



# Water Quality vs Functionality

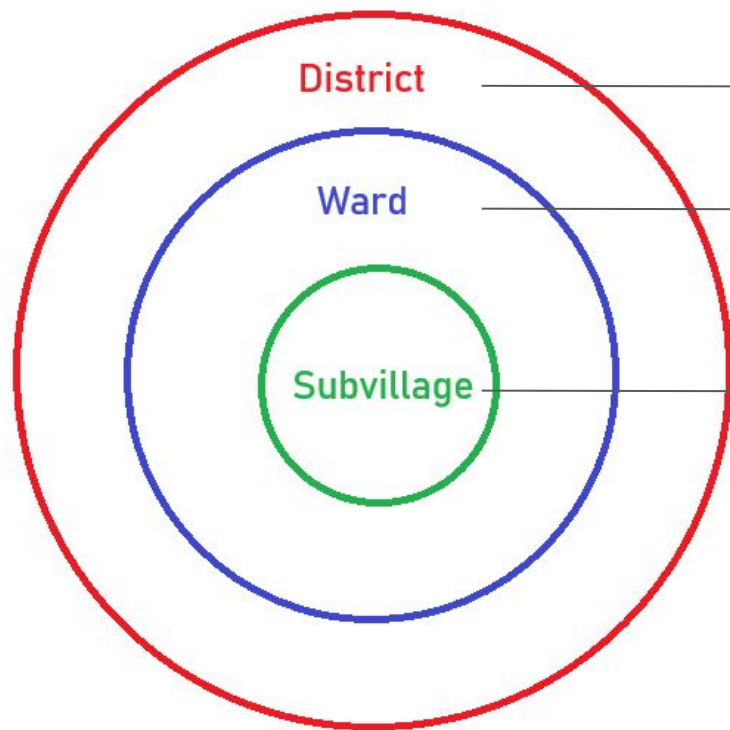


# Extraction Type





# “Group fill” missing values



	district_code	latitude	longitude
4	1	-1.825359	31.130847
9	1	-1.257051	30.626991
21	1	-1.541205	30.878919

		latitude	longitude
19118	bumera	-1.327684	34.280274
42088	bumera	-1.326142	34.276339
1575	bumera	-1.326913	34.278307

Check for NaNs,  
If NaNs present, move to next  
group (outer circle)

	subvillage	latitude	longitude
304	barazani	-4.389084	33.430917
5782	barazani	-4.405503	33.43444
31403	barazani	-4.397294	33.432679

# Columns selected for initial training

Col. No	Column Name Description	Processing
2	<u>date_recorded</u> - The date the row was entered	outliers removed (31 records) 2004 30 2002 1 Year extracted from date, Ordinal encoding is used
4	<u>gps_height</u> - Altitude of the well	Boxplot is used to check for outlier, No outlier MinMax normalization is used
6	longitude - GPS coordinate	Outlier removed using IQR formula and replaced with median <u>MinMax Normalization</u>
7	latitude - GPS coordinate	No outlier but there is a gap between -1 and 0 in <u>th</u> box plot (1819 values close to 0) <u>MinMax Normalization</u> is used
10	basin - Geographic water basin	09 unique values OneHotEncoding is used
12/13	<u>Region/region_code</u> - Geographic location (coded)	21/27 unique values I used one hot encoding

# Columns selected for initial training

Col. No	Column Name Description	Processing
14/15	<u>district_code</u> /lga - Geographic location	20/125 unique values I used one hot encoding
17	Population	Outliers replaced with mean
18	public_meeting - True/False	3334 NaN values replaced with string Unknown Now we have three categories,Used Onehot encoding
22	permit - If the waterpoint is permitted	3055 NaN values replaced with string Unknown Now we have three categories,Used Onehot encoding
23	construction_year - Year the waterpoint was constructed	33% values are zero, mean its unknown replaced with median of non zero values
24	extraction_type - The kind of extraction the waterpoint uses	<u>extraction type</u> ,group and class has 18,13 and 7 categories respectively, extraction type is selected,1 hot encoding is used

# Columns selected for initial training

Col. No	Column Name Description	Processing
27	management - How the <u>waterpoint</u> is managed	management has 12 categories, group has 5 categories management column is selected, 1 hot encoding is used
29	payment - What the water costs	Both payment and payment type have the same 7 unique values Payment column selected and, 1 hot encoding is used
31	<u>water_quality</u> - The quality of the water	water quality has 8 unique values, quality_group has 6 unique values water quality is selected, 1 hot encoding is used
33	quantity - The quantity of water	quantity and <u>quantity_group</u> columns are same with 5 categories I chose <u>quantity</u> , 1 hot encoding is used
35	source - The source of the water	source has 10 unique values, <u>source_type</u> has 7 unique values source is selected, 1 hot encoding is used
37	source_class - The source of the water	3 unique categories 1 hot encoding is used
38	waterpoint_type - The kind of waterpoint	water point type and group has 7 and 6 unique values, water point type is selected, 1 hot encoding is used

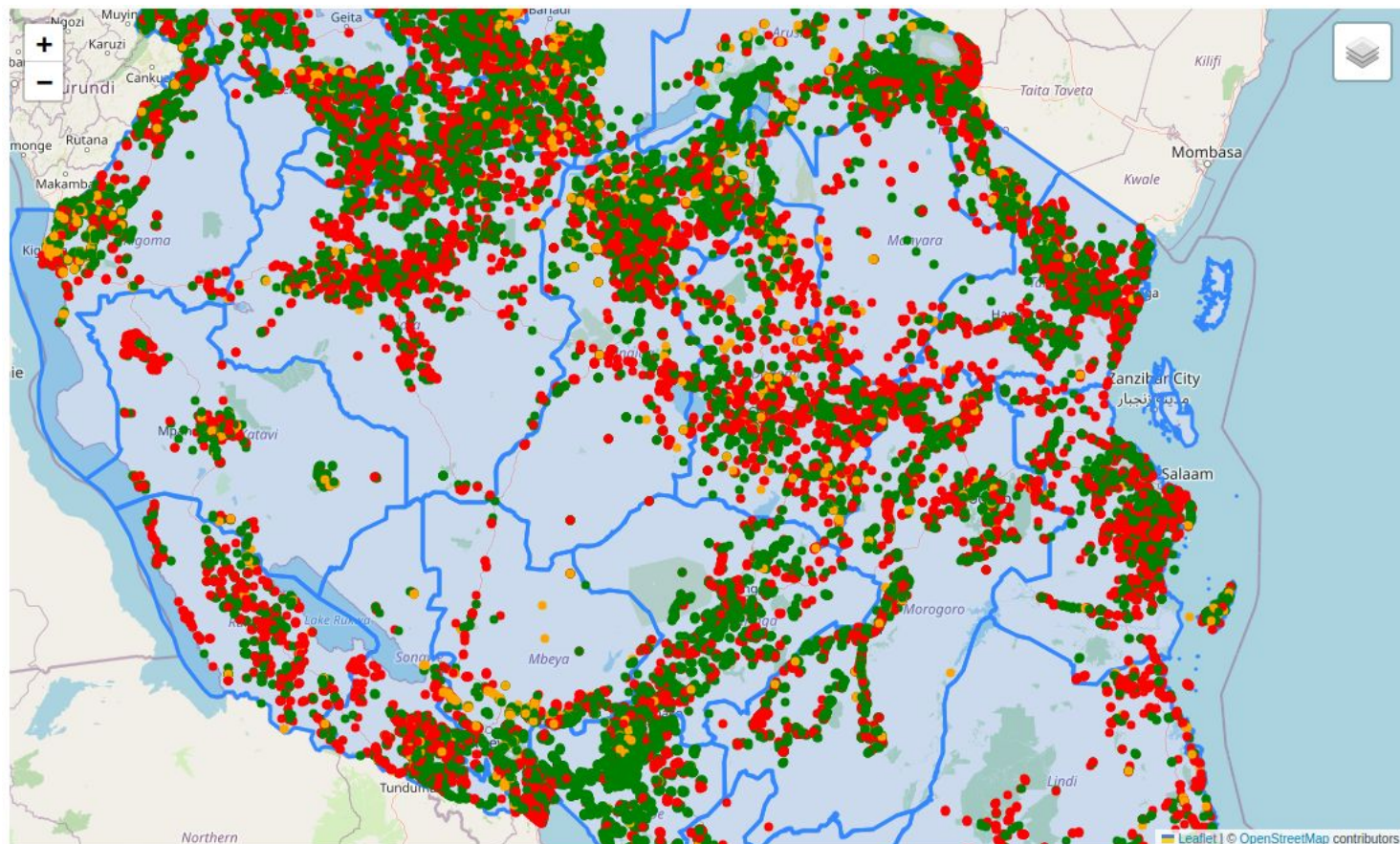
# Columns dropped for initial training

Col. No	Column Name Description	Reason
1	<u>amount_tsh</u> - Total static head (amount water available to <u>waterpoint</u> )	50 % values are zero
3	<u>funder</u> - Who funded the well	1895 Unique values.3535 <u>NaN</u> values. Highest contribution of pumps (around 8000) from <u>Govt of Tanzania</u>
5	<u>installer</u> - Organization that installed the well	2143 Unique values, 3653 <u>NaN</u> values. <u>Higest</u> contribution from DWE (around 17000)
8	<u>wpt_name</u> -water point nam	37381 unique values,2 Nan values
9	<u>num_private</u> -	70 % values are zero
11	<u>subvillage</u> - Geographic location	19281 unique values, 371 Nan values
16	<u>ward</u> - Geographic location	2092 unique values, No <u>NaN</u> values



# Columns dropped for initial training

Col. No	Column Name Description	Reason
19	<u>recorded_by</u> - Group entering this row of data	All values are same
20	<u>scheme_management</u> - Who operates the <u>waterpoint</u>	2695 unique values, 3874 <u>NaN</u> values
21	<u>scheme_name</u> - Who operates the <u>waterpoint</u>	50% values are <u>NaN</u> (count=28790)
25/26	<u>extraction_type_group</u> / <u>extraction_type_group</u> - The kind of extraction the <u>waterpoint</u> uses	13/7 categories..already selected <u>extraction_type</u> with 18 categories
28	<u>management_group</u> - How the <u>waterpoint</u> is managed	5 unique values in management group. Selected management column with 12 categories
30	<u>payment_type</u> - What the water costs	07 unique <u>values</u> ..almost same as selected payment column
32,34	Quality group/quantity group	water quality/quantity columns are selected
36,39	Source type, water point type group	Source/ <u>waterpoint_type</u> columns are selected



# Conclusion

Plenty of missing data that needs to be handled

Imbalance in data

Some features tend to have interesting relation to target variable