

# 1 Business Problem:

## 1.1 Background

In Tanzania Half of the population does not have access to water pumps through coordination of various charities, others require repair.

## 1.2 Problem Statement

We need to develop a ternary classification model to predict repairs, and which don't work at all. An understanding of how much potable water is available to Tanzanian communities.

When examining what causes a waterpoint to fail we

- What kind of pump is operating
- When it was installed
- Where it was installed
- Who installed it
- How it is managed

## 2 Import Libraries

### Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small categories
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import matplotlib.ticker as mtick
import seaborn as sns
import numpy as np
import scipy.stats as stats
import statsmodels.api as sm
# import xgboost as xgb
import catboost
import time
import warnings
warnings.filterwarnings('ignore')

from sklearn.utils import class_weight
from sklearn.metrics import accuracy_score
from catboost import Pool, sum_models
from catboost import CatBoostClassifier
from statsmodels.formula.api import ols
from sklearn.feature_selection import RF
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import KFold
from sklearn.preprocessing import LabelEncoder
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from sklearn import metrics
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import uniform, truncnorm
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import LabelEncoder
```

## 3 Data Exploration

```
In [2]: df_test = pd.read_csv('data/Test-set-validation')
```

```
In [3]: df_train_set = pd.read_csv('data/Trainin
df_train_labels = pd.read_csv('data/Trai
```

### 3.1 Data Fields

Predict one of the three classes based on a number c

The following set of information about waterpoints is

- amount\_tsh — Total static head (zero for open ta
- date\_recorded — The date the row was entered
- funder — Who funded the well
- gps\_height — Altitude of the well
- installer — Organization that installed the well
- longitude — GPS coordinate
- latitude — GPS coordinate
- wpt\_name — Name of the waterpoint if there is c
- num\_private — No information
- basin — Geographic water basin
- subvillage — Geographic location
- region — Geographic location
- region\_code — Geographic location (coded)
- district\_code — Geographic location (coded)
- lga — Geographic location
- ward — Geographic location
- population — Population around the well
- public\_meeting — True/False
- recorded\_by — Group entering this row of data
- scheme\_management — Who operates the wate
- scheme\_name — Who operates the waterpoint
- permit — If the waterpoint is permitted
- construction\_year — Year the waterpoint was co
- extraction\_type — The kind of extraction the wat
- extraction\_type\_group — The kind of extraction t
- extraction\_type\_class — The kind of extraction tl
- management — How the waterpoint is managed
- management\_group — How the waterpoint is ma
- payment — What the water costs
- payment\_type — What the water costs
- water\_quality — The quality of the water
- quality\_group — The quality of the water
- quantity — The quantity of water
- quantity\_group — The quantity of water (duplicat
- source — The source of the water
- source\_type — The source of the water
- source\_class — The source of the water
- waterpoint\_type — The kind of waterpoint
- waterpoint\_type\_group — The kind of waterpoint

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [4]: `df_train_set.head()`

Out[4]:

	id	amount_tsh	date_recorded	funder	gps_height
0	69572	6000.0	2011-03-14	Roman	
1	8776	0.0	2013-03-06	Grumeti	
2	34310	25.0	2013-02-25	Lottery Club	
3	67743	0.0	2013-01-28	Unicef	
4	19728	0.0	2011-07-13	Action In A	

5 rows × 6 columns

In [5]: `df_train_set.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 59400 entries, 0 to 59399
Data columns (total 40 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     59400 non-null  object
1   amount_tsh                           59400 non-null  float64
2   date_recorded                         59400 non-null  object
3   funder                                55765 non-null  object
4   gps_height                            59400 non-null  float64
5   installer                             55745 non-null  object
6   longitude                             59400 non-null  float64
7   latitude                              59400 non-null  float64
8   wpt_name                             59400 non-null  object
9   num_private                           59400 non-null  float64
10  basin                                 59400 non-null  object
11  subvillage                            59029 non-null  object
12  region                                59400 non-null  object
13  region_code                           59400 non-null  object
14  distance_to_nearest_well              59400 non-null  float64
```

In [6]: `df_train_labels['status_group'].value_counts()`

Out[6]:

functional	0.543081
non functional	0.384242
functional needs repair	0.072677

Name: status\_group, dtype: float64

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [7]: df\_train\_labels.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 59400 entries, 0 to 59399
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  -
0   id              59400 non-null   int64
1   status_group    59400 non-null   object
dtypes: int64(1), object(1)
memory usage: 928.2+ KB
```

In [8]: df\_train\_labels['status\_group'].value\_counts()

```
Out[8]: functional          0.543081
non functional             0.384242
functional needs repair    0.072677
Name: status_group, dtype: float64
```

- Most wells are either functional or non functional

In [9]: # print top 5 most frequent values in each column

```
for col in df_train_set.columns:
    print(col, '\n', df_train_set[col].value_counts().head(5))
```

```
id
2047      0.000017
72310     0.000017
49805     0.000017
51852     0.000017
62091     0.000017
Name: id, dtype: float64
```

```
amount_tsh
0.0        0.700993
500.0      0.052222
50.0       0.041616
1000.0     0.025051
20.0       0.024630
Name: amount_tsh, dtype: float64
```

```
date_recorded
2011-02-15      0.000000
```

- amount\_tsh: majority are zeros (70%)
- gps\_height: top value is zero (34%)
- longitude: top value is zero (3%) these are not valid
- region and region\_code: remove region
- wpt\_name: top value is none. (5%)
- num\_private: top value is zero (98%)

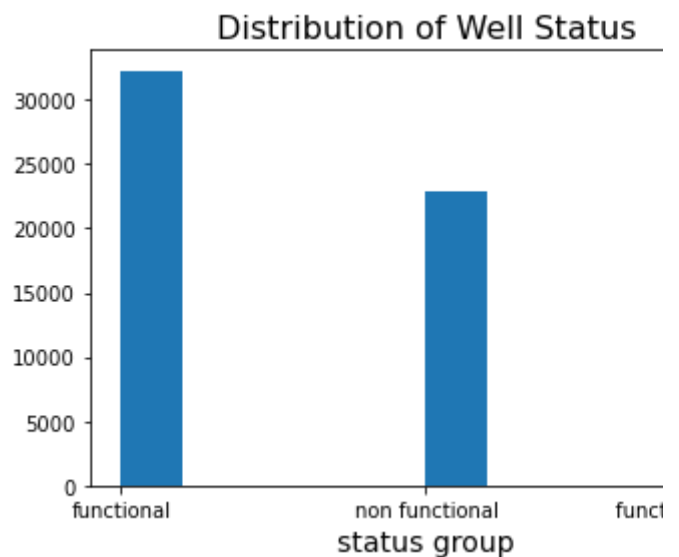
- population: top value is zero (36%)
- recorded\_by: only one value can remove
- construction\_year: top value is zero (34%)
- extraction\_type and extraction\_type\_field are the
- quality\_group looks to be a replacement for water
- quantity and quantity\_group look the same
- source vs source\_type vs source\_class: source\_class
- waterpoint\_type vs waterpoint\_type\_group: waterpoint\_type

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [10]: # Plot of the target status group

```
plt.hist(df_train_labels['status_group'])
plt.xlabel('status group', fontsize=14)
plt.title("Distribution of Well Status",
plt.show()
```



Here we see a strong class imbalance that will need to

## 4 Data Preparation

### 4.1 Missing Values

#### 4.1.1 permit

In [11]: df\_train\_set['permit'].value\_counts(normalized=True)

```
Out[11]: True      0.654074
False    0.294478
NaN       0.051448
Name: permit, dtype: float64
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [12]: # df_train_set.subvillage.fillna(0, inplace=True)
df_train_set.permit.fillna(0, inplace=True)
# df_train_set.public_meeting.fillna(0, inplace=True)
# df_train_set.scheme_management.fillna(0, inplace=True)
```

```
In [13]: df_train_set['permit'].value_counts(normalize=True)
```

```
Out[13]: True      0.654074
False     0.345926
Name: permit, dtype: float64
```

```
In [14]: #Make Permit boolean
df_train_set['permit'] = df_train_set['permit'].astype(bool)
df_train_set['permit'].value_counts(normalize=True)
```

```
Out[14]: 1      0.654074
0      0.345926
Name: permit, dtype: float64
```

### 4.1.2 construction\_year

```
In [15]: # Close to 35% are zero lets replace th
df_train_set['construction_year'].value_
```

```
Out[15]: 0      0.348636
          2010     0.044529
          2008     0.043990
          2009     0.042643
          2000     0.035202
          2007     0.026717
          2006     0.024764
          2003     0.021650
          2011     0.021145
          2004     0.018906
          2012     0.018249
          2002     0.018098
          1978     0.017458
          1995     0.017071
          2005     0.017020
          1999     0.016481
          1998     0.016263
          1990     0.016061
          1985     0.015909
          1980     0.013653
          1996     0.013653
          1984     0.013114
          1982     0.012525
          1994     0.012424
          1972     0.011919
          1974     0.011380
          1997     0.010842
          1992     0.010774
          1993     0.010236
          2001     0.009091
          1988     0.008771
          1983     0.008215
          1975     0.007357
          1986     0.007306
          1976     0.006970
          1970     0.006919
          1991     0.005455
          1989     0.005320
          1987     0.005084
          1981     0.004007
          1977     0.003401
          1979     0.003232
          1973     0.003098
          2013     0.002963
          1971     0.002441
          1960     0.001717
          1967     0.001481
          1963     0.001431
          1968     0.001296
          1969     0.000993
          1964     0.000673
          1962     0.000505
          1961     0.000354
          1965     0.000320
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



1966 0.000286

Name: construction\_year, dtype: float64

In [17]: nonzero\_median\_year = df\_train\_set[ df\_t

In [18]: nonzero\_median\_year['construction\_year']  
nonzero\_median\_year['construction\_year']

Out[18]: 2000.0

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [20]: `df_train_set['construction_year'].value_`

Out[20]:

0	0.348636
2010	0.044529
2008	0.043990
2009	0.042643
2000	0.035202
2007	0.026717
2006	0.024764
2003	0.021650
2011	0.021145
2004	0.018906
2012	0.018249
2002	0.018098
1978	0.017458
1995	0.017071
2005	0.017020
1999	0.016481
1998	0.016263
1990	0.016061
1985	0.015909
1980	0.013653
1996	0.013653
1984	0.013114
1982	0.012525
1994	0.012424
1972	0.011919
1974	0.011380
1997	0.010842
1992	0.010774
1993	0.010236
2001	0.009091
1988	0.008771
1983	0.008215
1975	0.007357
1986	0.007306
1976	0.006970
1970	0.006919
1991	0.005455
1989	0.005320
1987	0.005084
1981	0.004007
1977	0.003401
1979	0.003232
1973	0.003098
2013	0.002963
1971	0.002441
1960	0.001717
1967	0.001481
1963	0.001431
1968	0.001296
1969	0.000993
1964	0.000673
1962	0.000505
1961	0.000354
1965	0.000320

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

1966 0.000286

Name: construction\_year, dtype: float64

## 4.2 Replace misspellings and group

### 4.2.1 Reformat Installer col

```
In [21]: df_train_set['installer'] = df_train_set
df_train_set['installer'].replace(to_repl
```

```
In [22]: df_train_set['installer'].value_counts(r
```

```
Out[22]: dwe                0.293013
other                0.074663
government          0.031835
hesawa              0.023485
rwe                 0.020303
...
kkkt ndrumanengi   0.000017
safari roya        0.000017
tsrc               0.000017
magul              0.000017
hemed abdallah     0.000017
Name: installer, Length: 1934, dtype: fl
```

```
In [23]: df_clean=pd.read_csv('data/lookups.csv')
df_train_set = pd.merge(df_train_set, df
df_train_set.head()
```

```
Out[23]:
```

	id	amount_tsh	date_recorded	funder	gps_t
0	69572	6000.0	2011-03-14	Roman	
1	8776	0.0	2013-03-06	Grumeti	
2	34310	25.0	2013-02-25	Lottery Club	
3	67743	0.0	2013-01-28	Unicef	
4	19728	0.0	2011-07-13	Action In A	

5 rows × 42 columns

```
In [24]: df_train_set['installer_new'].isna().sun
```

```
Out[24]: 44958
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [25]: df_train_set['installer'].value_counts()
```

```
Out[25]: dwe          0.293013
other         0.074663
government    0.031835
hesawa        0.023485
rwe           0.020303
...
kkkt ndrumangeni 0.000017
safari roya     0.000017
tsrc            0.000017
magul           0.000017
hemed abdallah  0.000017
Name: installer, Length: 1934, dtype: float64
```

```
In [26]: df_train_set['installer_new'].fillna(df_train_set
```

```
Out[26]:
```

	id	amount_tsh	date_recorded	funder
0	69572	6000.0	2011-03-14	Roman
1	8776	0.0	2013-03-06	Grumeti
2	34310	25.0	2013-02-25	Lottery Club
3	67743	0.0	2013-01-28	Unicef
4	19728	0.0	2011-07-13	Action In A
...	...	...	...	...
59395	60739	10.0	2013-05-03	Germany Republi
59396	27263	4700.0	2011-05-07	Cefa-njombe
59397	37057	0.0	2011-04-11	NaN
59398	31282	0.0	2011-03-08	Malec
59399	26348	0.0	2011-03-23	World Bank

59400 rows x 42 columns

```
In [27]: df_train_set['installer_new'].isna().sum
```

```
Out[27]: 0
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [28]: df_train_set['installer_new'].value_counts()
```

```
Out[28]: dwe      0.304983
          other    0.074663
          tanzanian government 0.062323
          community 0.033266
          danida   0.028300
          ...
          b.a.p    0.000017
          lindi contractor 0.000017
          ji       0.000017
          friedkin conservation fund 0.000017
          hemed abdallah 0.000017
          Name: installer_new, Length: 1652, dtype: object
```

```
In [29]: del df_train_set['installer']
          del df_train_set['installer_old']
```

```
In [30]: df_installer_cnt = df_train_set['installer_new'].value_counts()
          df_installer_cnt
```

```
Out[30]: dwe      0.304983
          other    0.074663
          tanzanian government 0.062323
          community 0.033266
          danida   0.028300
          ...
          b.a.p    0.000017
          lindi contractor 0.000017
          ji       0.000017
          friedkin conservation fund 0.000017
          hemed abdallah 0.000017
          Name: installer_new, Length: 1652, dtype: object
```

```
In [31]: df_installer_cnt.head(20)
```

```
Out[31]: dwe      0.304983
          other    0.074663
          tanzanian government 0.062323
          community 0.033266
          danida   0.028300
          hesawa   0.023620
          council  0.022929
          rwe      0.020303
          church   0.016330
          kkkt     0.015623
          finw     0.013199
          tcrs     0.012458
          world vision 0.012037
          ces      0.010269
          amref    0.007458
          lga      0.006953
          tasaf    0.006919
          wedeco   0.006700
          dmdd     0.006330
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [32]: other_list = df_installer_cnt[df_installer_cnt['other'] > 0]
other_list
other
'norad',
'unicef',
'twesa',
'da',
'wu',
'acra',
'sema',
'jaica',
'oxfam',
'shipo',
'local',
'idara ya maji',
'villagers',
'sengerema water department',
'kiliwater',
'dh',
'kuwait',
'distri',
'lawatefuka water sup',
'magadini-makiwaru wa',
'...
```

```
In [33]: df_train_set['installer_new'].replace(to_replace=df_train_set['installer_new'].value_counts().index, value=0)
```

```
Out[33]: other          0.424360
dwe                    0.304983
tanzanian government  0.062323
community              0.033266
danida                 0.028300
hesawa                 0.023620
council                0.022929
rwe                    0.020303
church                 0.016330
kkkt                   0.015623
finw                   0.013199
tcrs                   0.012458
world vision           0.012037
ces                    0.010269
Name: installer_new, dtype: float64
```

In [34]: `df_train_set.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 59400 entries, 0 to 59399
Data columns (total 40 columns):
#   Column                                Non-Null Count
---  --
0   id                                    59400 non-null
1   amount_tsh                          59400 non-null
2   date_recorded                        59400 non-null
3   funder                               55765 non-null
4   gps_height                           59400 non-null
5   longitude                            59400 non-null
6   latitude                             59400 non-null
7   wpt_name                             59400 non-null
8   num_private                          59400 non-null
9   basin                                59400 non-null
10  subvillage                           59029 non-null
11  region                                59400 non-null
12  region_code                           59400 non-null
13  district_code                         59400 non-null
14  lga                                    59400 non-null
15  ward                                  59400 non-null
16  population                            59400 non-null
17  public_meeting                        56066 non-null
18  recorded_by                           59400 non-null
19  scheme_management                     55523 non-null
20  scheme_name                           31234 non-null
21  permit                                59400 non-null
22  construction_year                     59400 non-null
23  extraction_type                       59400 non-null
24  extraction_type_group                  59400 non-null
25  extraction_type_class                   59400 non-null
26  management                             59400 non-null
27  management_group                       59400 non-null
28  payment                                59400 non-null
29  payment_type                           59400 non-null
30  water_quality                          59400 non-null
31  quality_group                          59400 non-null
32  quantity                               59400 non-null
33  quantity_group                         59400 non-null
34  source                                 59400 non-null
35  source_type                            59400 non-null
36  source_class                           59400 non-null
37  waterpoint_type                        59400 non-null
38  waterpoint_type_group                  59400 non-null
39  installer_new                          59400 non-null
dtypes: float64(3), int64(8), object(29)
memory usage: 18.6+ MB
```

## 4.2.2 Reformat Funder col

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [35]: df_train_set['funder'] = df_train_set['funder'].replace(to_replace='funder', value='government of tanzania')
```

```
In [36]: df_train_set['funder'].value_counts(normalize=True)
```

```
Out[36]: government of tanzania    0.152929
other                             0.074276
danida                           0.052424
hesawa                           0.037071
rwssp                             0.023131
...
tlc/emmanuel kasoga              0.000017
villlage contributi              0.000017
hasnan murig (mbunge)           0.000017
samweli mshosha                 0.000017
kajima                          0.000017
Name: funder, Length: 1897, dtype: float64
```

```
In [37]: df_clean_funder=pd.read_csv('data/lookug')
df_train_set = pd.merge(df_train_set, df_clean_funder, on='id')
df_train_set.head()
```

```
Out[37]:
```

	id	amount_tsh	date_recorded	funder	gps_h
0	69572	6000.0	2011-03-14	roman	
1	8776	0.0	2013-03-06	grumeti	
2	34310	25.0	2013-02-25	lottery club	
3	67743	0.0	2013-01-28	unicef	
4	19728	0.0	2011-07-13	action in a	

5 rows × 42 columns

```
In [38]: df_train_set['funder_new'].isna().sum()
```

```
Out[38]: 39498
```



## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [39]: `df_train_set['funder'].value_counts(normalized=True)`

```
Out[39]: government of tanzania    0.152929
other                            0.074276
danida                          0.052424
hesawa                          0.037071
rwssp                           0.023131
...
tlc/emmanuel kasoga             0.000017
villlage contributi            0.000017
hasnan murig (mbunge)          0.000017
samweli mshosha                0.000017
kajima                         0.000017
Name: funder, Length: 1897, dtype: float64
```

In [40]: `df_train_set['funder_new'].fillna(df_train_set['funder']).value_counts(normalized=True)`

```
Out[40]:
```

	id	amount_tsh	date_recorded	funder
0	69572	6000.0	2011-03-14	roman
1	8776	0.0	2013-03-06	grumeti
2	34310	25.0	2013-02-25	lottery club
3	67743	0.0	2013-01-28	unicef
4	19728	0.0	2011-07-13	action in a
...	...	...	...	...
59395	60739	10.0	2013-05-03	germany republi
59396	27263	4700.0	2011-05-07	cefa-njombe
59397	37057	0.0	2011-04-11	other
59398	31282	0.0	2011-03-08	malec
59399	26348	0.0	2011-03-23	world bank

59400 rows × 42 columns

```
In [41]: df_train_set['funder_new'].isna().sum()
```

```
Out[41]: 0
```

```
In [42]: df_train_set['funder_new'].value_counts()
```

```
Out[42]: government      0.156869
other                  0.074276
danida                 0.052559
hesawa                 0.037323
kkkt                   0.026027
...
chongolo               0.000017
lgdbg                  0.000017
petro patrice          0.000017
mwakalinga             0.000017
quicklw                0.000017
Name: funder_new, Length: 1683, dtype: float64
```

```
In [43]: del df_train_set['funder']
del df_train_set['funder_old']
```

```
In [44]: df_installer_cnt_funder = df_train_set['funder_new']
```

```
Out[44]: government      0.156869
other                  0.074276
danida                 0.052559
hesawa                 0.037323
kkkt                   0.026027
...
chongolo               0.000017
lgdbg                  0.000017
petro patrice          0.000017
mwakalinga             0.000017
quicklw                0.000017
Name: funder_new, Length: 1683, dtype: float64
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [45]: df_installer_cnt_funder.head(20)
```

```
Out[45]: government      0.156869
          other          0.074276
          danida         0.052559
          hesawa         0.037323
          kkkt           0.026027
          district council 0.025051
          church         0.024529
          rwssp          0.023131
          world bank     0.022710
          world vision   0.021498
          unicef         0.019798
          tasaf          0.014764
          dhv            0.013956
          private individual 0.013906
          dwsp           0.013653
          norad          0.012879
          finw           0.012559
          german group    0.012492
          tcrs           0.010455
          ministry of water 0.009933
          Name: funder_new, dtype: float64
```

```
In [46]: other_list_funder = df_installer_cnt_funder['funder_new']
          other_list_funder
```

```

          'halmashaur',
          'br',
          'professor ben ohio university',
          'member of parliament',
          'i.e.c',
          'lawate fuka water suppl',
          'abd',
          'sauwasa',
          'mws',
          'school',
          'parastatal',
          'menon',
          'summit for water',
          'rwssp',
          'local',
          'unicef/ csp',
          'aficare',
          'hapa',
          'el',
          ...
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [47]: df_train_set['funder_new'].replace(to_replace=df_train_set['funder_new'].value_counts(
```

```
Out[47]: other                0.485842
government                0.156869
danida                    0.052559
hesawa                    0.037323
kkkt                      0.026027
district council          0.025051
church                    0.024529
rwssp                     0.023131
world bank                 0.022710
world vision              0.021498
unicef                    0.019798
tasaf                     0.014764
dhv                       0.013956
private individual        0.013906
dwsp                      0.013653
norad                     0.012879
finw                      0.012559
german group              0.012492
tcrs                      0.010455
Name: funder_new, dtype: float64
```

### 4.2.3 Group Other Remaining Columns

```
In [48]: def group_other(col,perc):
df_train_set[col] = df_train_set[col]
df_cnt = df_train_set[col].value_counts()
other_list = df_cnt[df_cnt<perc].index

return df_train_set[col].replace(to_replace=other_list, value='other')
print(df_train_set[col].value_counts())
```

#### 4.2.3.1 lga

```
In [49]: df_train_set['lga'].value_counts(normalized=True)
```

```
Out[49]: Njombe                0.042138
Arusha Rural                0.021077
Moshi Rural                 0.021061
Bariadi                    0.019815
Rungwe                     0.018620
...
Moshi Urban                0.001330
Kigoma Urban               0.001195
Arusha Urban               0.001061
Lindi Urban                0.000354
Nyamagana                  0.000017
Name: lga, Length: 125, dtype: float64
```

```
In [50]: group_other('lga',.01)
```

```
In [51]: df_train_set['lga'].value_counts(normali
```

```
Out[51]: other          0.475438
njombe          0.042138
arusha rural    0.021077
moshi rural     0.021061
bariadi         0.019815
rungwe          0.018620
kilosa          0.018418
kasulu          0.017626
mbozi           0.017407
meru            0.016987
bagamoyo        0.016785
singida rural   0.016751
kilombero       0.016145
same            0.014764
kibondo         0.014714
kyela           0.014461
kahama          0.014074
magu            0.013872
kigoma rural    0.013872
maswa           0.013620
karagwe         0.012980
mbinga          0.012626
iringa rural    0.012256
serengeti       0.012054
lushoto         0.011684
namtumbo        0.011684
songea rural    0.011667
mpanda          0.011431
mvomero         0.011296
ngara           0.011263
ulanga          0.011195
makete          0.010606
kwimba          0.010556
mbarali         0.010539
hai             0.010522
rombo           0.010000
Name: lga, dtype: float64
```

## 4.3 Columns to drop

### 4.3.1 Mostly Empty

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [52]: df_train_set['num_private'].value_counts
```

```
Out[52]: 0          0.987256
          6          0.001364
          1          0.001229
          5          0.000774
          8          0.000774
          ...
         180         0.000017
         213         0.000017
          23         0.000017
          55         0.000017
          94         0.000017
Name: num_private, Length: 65, dtype: float64
```

### 4.3.2 Many Individual Values

```
In [53]: df_train_set['wpt_name'].value_counts(normalized=True)
```

```
Out[53]: none          0.059983
          Shuleni       0.029428
          Zahanati      0.013973
          Msikitini     0.009007
          Kanisani      0.005438
          ...
          Kwa Mzee Zenubius 0.000017
          Msanya        0.000017
          Kwa Mbumbuli   0.000017
          Irkisale Secondary School 0.000017
          Kwa Kiambae    0.000017
Name: wpt_name, Length: 37400, dtype: object
```

```
In [54]: df_train_set['ward'].value_counts(normalized=True)
```

```
Out[54]: Igosi         0.005168
          Imalinyi      0.004242
          Siha Kati     0.003906
          Mdandu        0.003889
          Nduruma       0.003653
          ...
          Uwanja wa Ndege 0.000017
          Nsemulwa       0.000017
          Kitete        0.000017
          Linda         0.000017
          Chinugulu     0.000017
Name: ward, Length: 2092, dtype: float64
```

### 4.3.3 Not Significant

**4.3.3.1 The features scheme\_management, quant\_date\_recorded, and recorded\_by will be deleted for**

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [55]: def drop(df):
# Columns to drop
df.drop([
    'scheme_management', 'quantity_c',
    'extraction_type', 'extraction_t',
    'subvillage', 'public_meeting',
    axis=1, inplace=True)
return df
```

To do: look at feat importance of different geographic

- basin — Geographic water basin
- subvillage — Geographic location
- region\_code — Geographic location (coded)
- district\_code — Geographic location (coded)
- lga — Geographic location
- ward — Geographic location

## 4.4 Categorical and Numerical

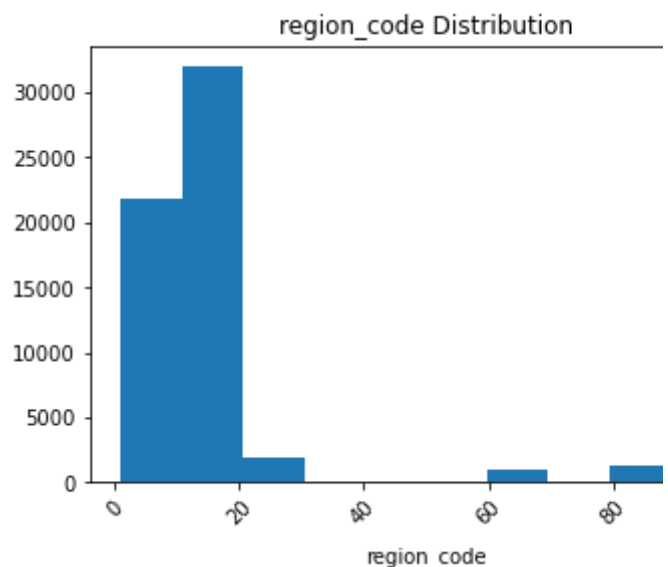
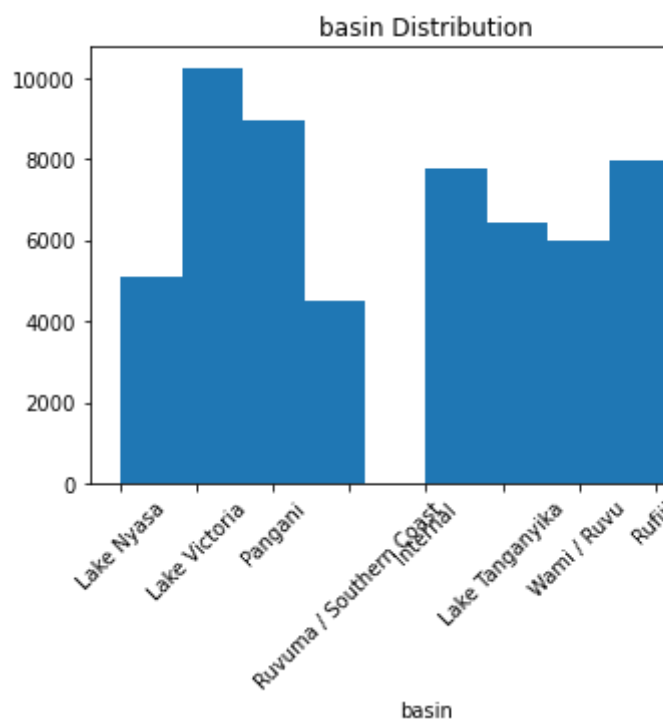
```
In [56]: categorical = ['basin', 'region_code', 'c',
                        'extraction_type_group', 'c',
                        'quality_group', 'quantit',
                        'installer_new', 'funder_
```



```
In [57]: for column in categorical:
plt.hist(df_train_set[column])
plt.xlabel(column)
plt.xticks(rotation=45)
plt.title("{} Distribution".format(c
plt.show()
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

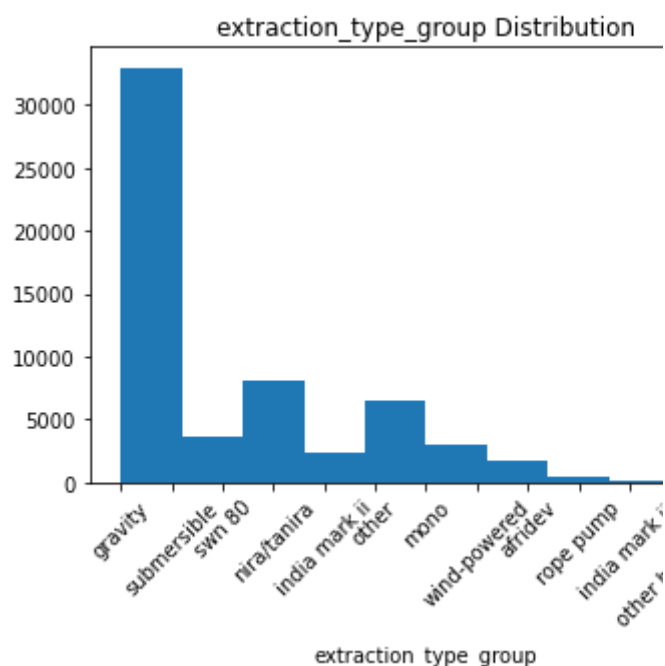
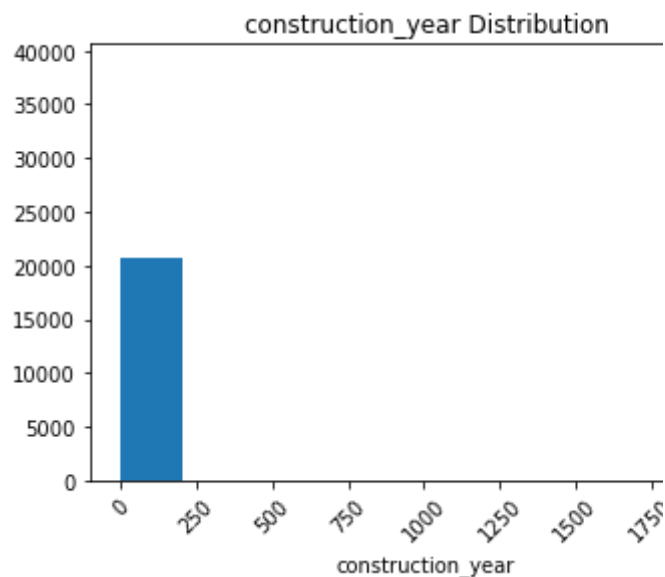






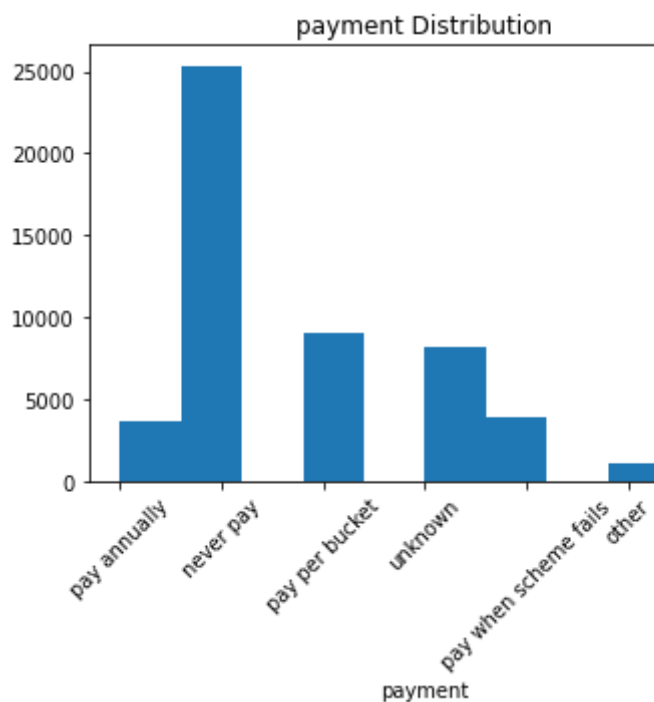
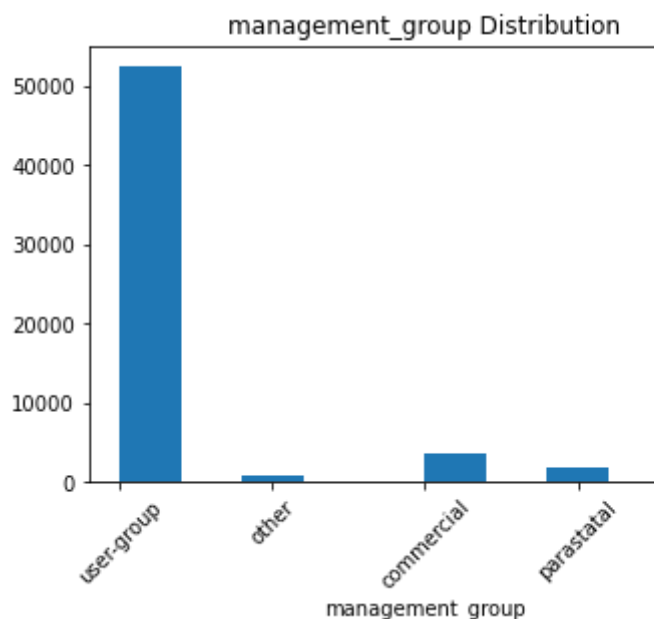
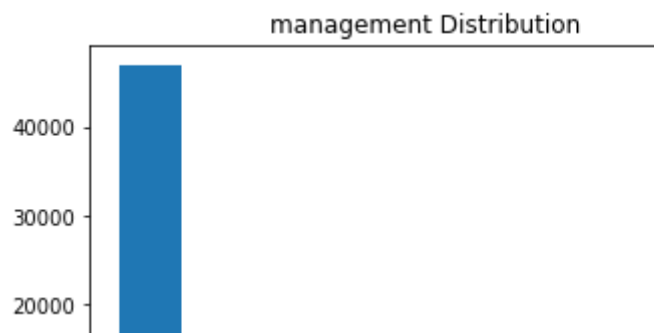
## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- 3 Data Exploration
  - 3.1 Data Fields
- 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



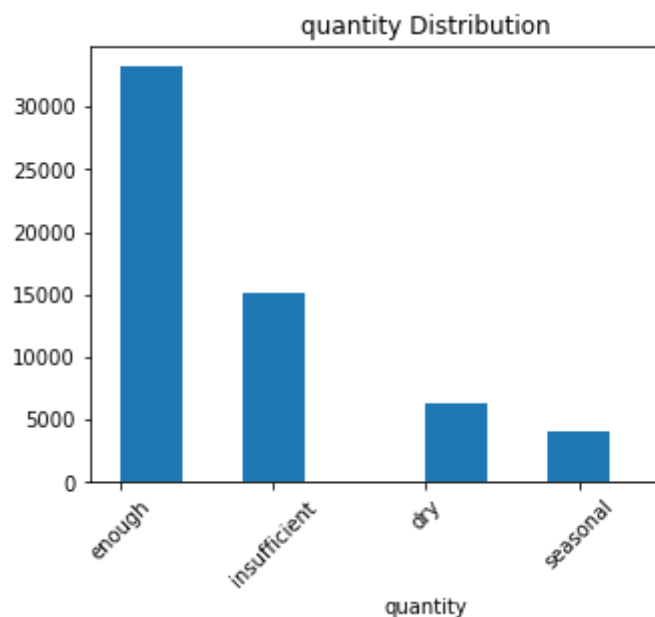
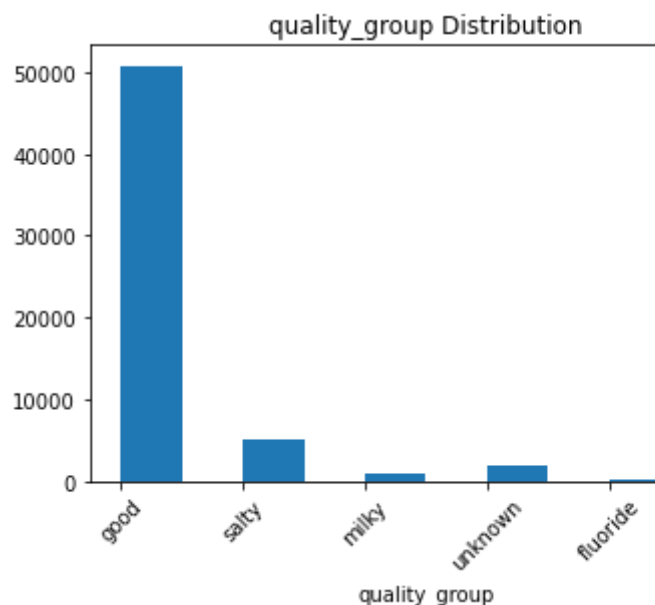
## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



## Contents 🔄 ⚙️

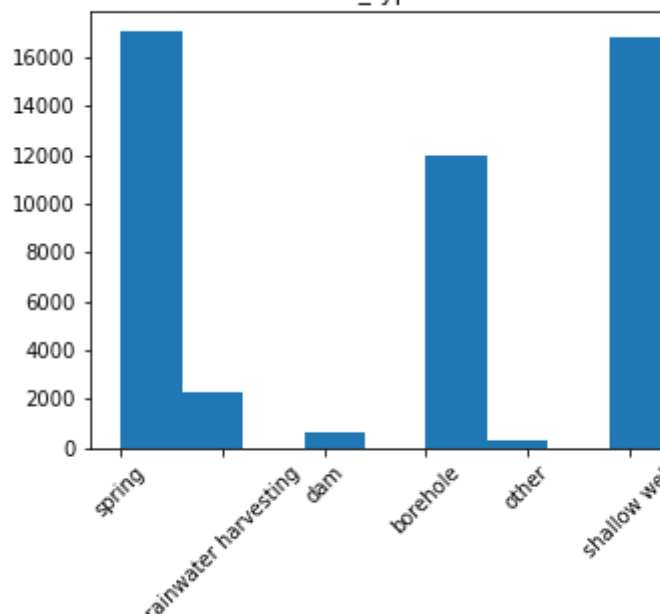
- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



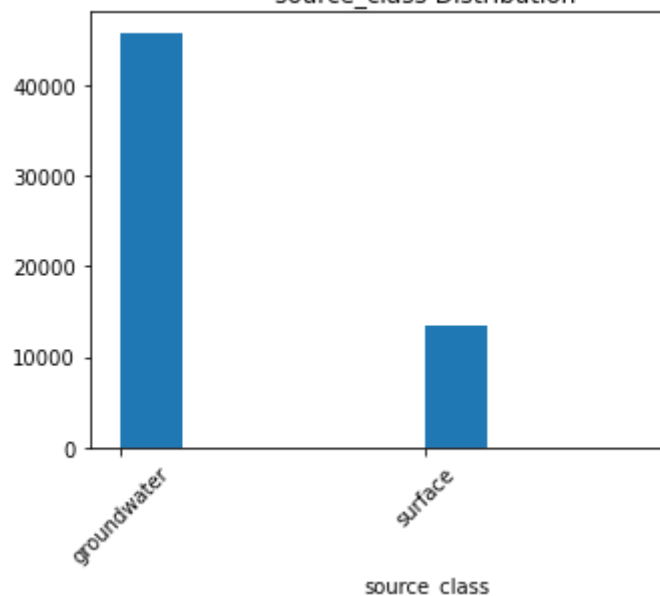
## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

source\_type Distribution

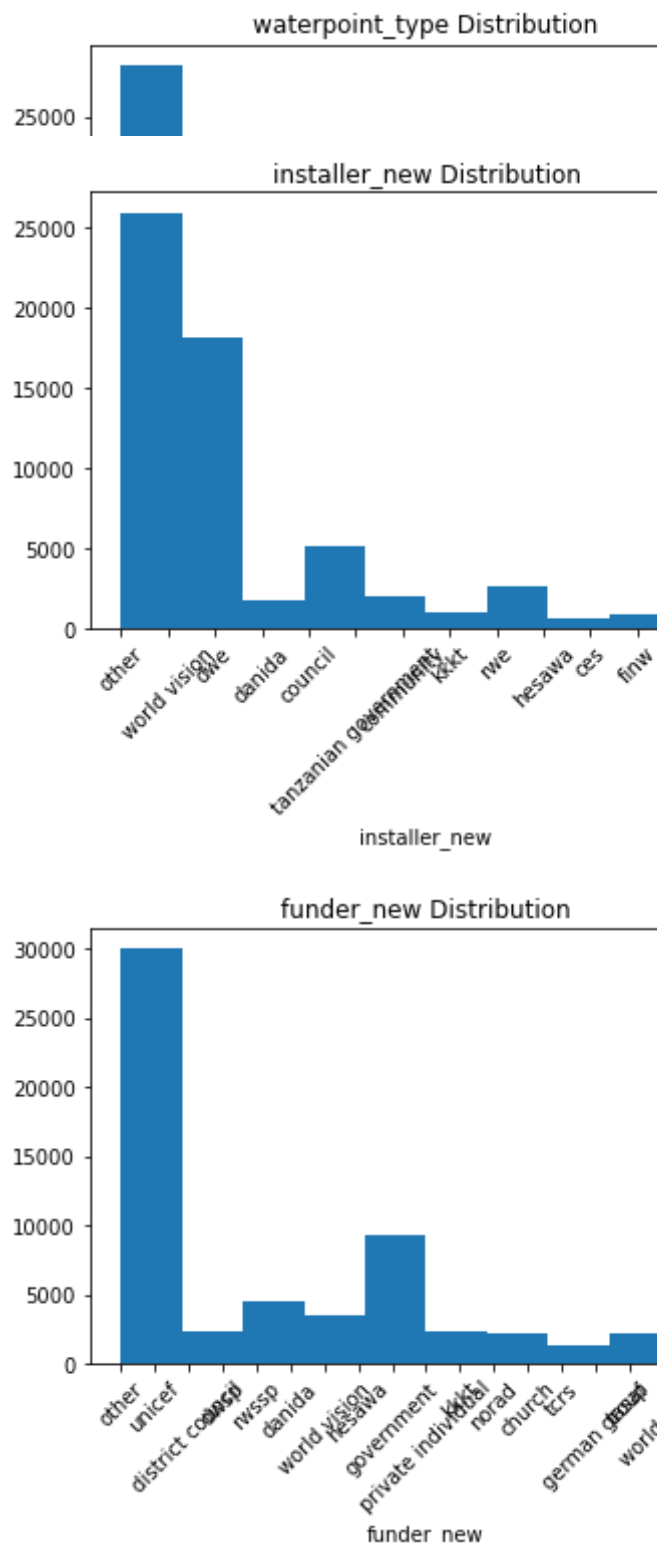


source\_class Distribution



## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning



### 4.4.1 Join Target: df\_train\_set to df\_train

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- 3 Data Exploration
  - 3.1 Data Fields
- 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [58]: df_train = df_train_set.join(df_train_large)
df_train = drop(df_train)

df_test = drop(df_test)
```

```
In [59]: df_train.info()
7    district_code      59400 non-na
8    lga                59400 non-na
9    population         59400 non-na
10   permit            59400 non-na
11   construction_year  59400 non-na
12   extraction_type_group  59400 non-na
13   management        59400 non-na
14   management_group   59400 non-na
15   payment            59400 non-na
16   quality_group      59400 non-na
17   quantity           59400 non-na
18   source_type        59400 non-na
19   source_class       59400 non-na
20   waterpoint_type     59400 non-na
21   installer_new      59400 non-na
22   funder_new         59400 non-na
23   id_right           47492 non-na
24   status_group       47492 non-na
dtypes: float64(4), int64(7), object(14)
memory usage: 11.8+ MB
```

```
In [60]: df_train.columns
```

```
Out[60]: Index(['id_left', 'amount_tsh', 'gps_height', 'region_code', 'district_code', 'construction_year', 'extraction_type_group', 'management_group', 'payment', 'source_type', 'source_class', 'waterpoint_type', 'installer_new', 'funder_new', 'id_right', 'status_group'],
              dtype='object')
```

```
In [61]: df_train['status_group'].value_counts(normalize=True)

Out[61]: functional                0.436229
non functional                    0.305522
NaN                               0.200471
functional needs repair           0.057778
Name: status_group, dtype: float64
```

```
In [62]: numerical = ['amount_tsh', 'gps_height', 'region_code', 'district_code', 'construction_year', 'extraction_type_group', 'management_group', 'payment', 'source_type', 'source_class', 'waterpoint_type', 'installer_new', 'funder_new', 'id_right', 'status_group']
```

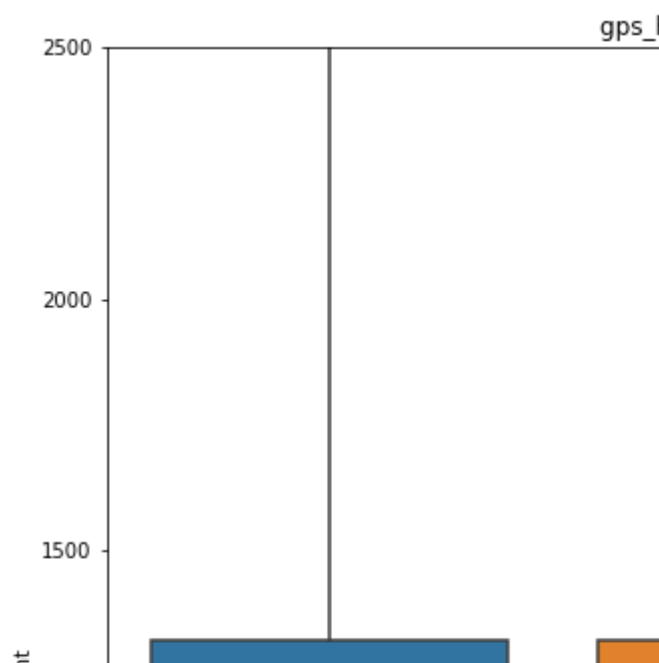
```
In [63]: numerical_large = ['gps_height', 'population', 'permit', 'construction_year', 'extraction_type_group', 'management_group', 'payment', 'source_type', 'source_class', 'waterpoint_type', 'installer_new', 'funder_new', 'id_right', 'status_group']
```

```
In [64]: numerical_small = ['amount_tsh']
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [65]: for variable in numerical_large:
          ax, figure = plt.subplots(1,1,figsize=(10,10))
          plt.ylim(-100,2500)
          sns.boxplot(x='status_group', y=variable, data=df_train)
          plt.title("{} vs. Well Condition".format(variable))
```



```
In [66]: del df_train['gps_height']
          del df_train['population']
```

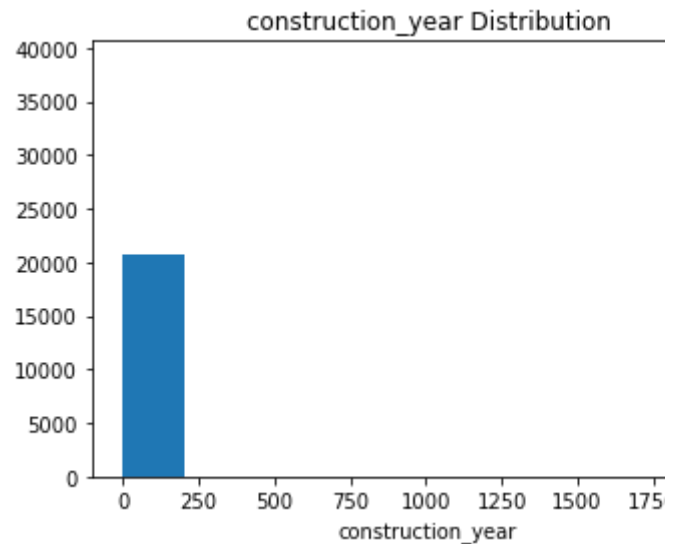
```
In [67]: for variable in numerical_small:
          ax, figure = plt.subplots(1,1,figsize=(10,10))
          plt.ylim(-1,200)
          sns.boxplot(x='status_group', y=variable, data=df_train)
          plt.title("{} vs. Well Condition".format(variable))
```



## 4.4.2 Column Binning

```
In [68]: year_plot = ['construction_year']
```

```
In [69]: for column in year_plot:
          plt.hist(df_train[column])
          plt.xlabel(column)
          plt.title("{} Distribution".format(c
          plt.show()
```



```
In [70]: bin_features = ['amount_tsh', 'constructi
```

```
In [71]: df_train['construction_year'].value_cour
```

```
Out[71]: 0      0.348636
          2010    0.044529
          2008    0.043990
          2009    0.042643
          2000    0.035202
          2007    0.026717
          2006    0.024764
          2003    0.021650
          2011    0.021145
          2004    0.018906
          2012    0.018249
          2002    0.018098
          1978    0.017458
          1995    0.017071
          2005    0.017020
          1999    0.016481
          1998    0.016263
          1990    0.016061
          1985    0.015909
          1988    0.015855
```

### Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [72]: df_train['amount_tsh'].value_counts(normalized=True)
```

```
Out[72]: 0.0          0.700993
          500.0        0.052222
          50.0         0.041616
          1000.0       0.025051
          20.0         0.024630
          ...
          8500.0       0.000017
          6300.0       0.000017
          220.0        0.000017
          138000.0     0.000017
          12.0         0.000017
Name: amount_tsh, Length: 98, dtype: float64
```

```
In [73]: df_train.amount_tsh.fillna(0, inplace=True)
```

```
In [74]: # and we will also take a look at their distribution
for feat in bin_features:
    print(df_train[feat].describe())
```

```
count      59400.000000
mean         317.650385
std         2997.574558
min           0.000000
25%           0.000000
50%           0.000000
75%          20.000000
max        350000.000000
Name: amount_tsh, dtype: float64
count      59400.000000
mean        1300.652475
std          951.620547
min           0.000000
25%           0.000000
50%          1986.000000
75%          2004.000000
max          2013.000000
Name: construction_year, dtype: float64
```

```
In [75]: # based on some intuition and the quantiles
tsh_bins = [-1, 0, 20, 350000]
construction_year_bins = [0, 1, 1985, 2005, 2013]
```

```
In [76]: # now lets convert these columns to binned
# note that when binning these, the default bins are 100
```

```
df_train['amount_tsh'] = pd.cut(df_train['amount_tsh'], tsh_bins)
df_train['construction_year'] = pd.cut(df_train['construction_year'], construction_year_bins)
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [77]: # and we will also take a look at the distribution
for feat in bin_features:
    print(df_train[feat].value_counts(normalized=True))

(-1, 0]          0.700993
(20, 350000]     0.249495
(0, 20]          0.049512
Name: amount_tsh, dtype: float64
NaN              0.348636
(1985.0, 2005.0] 0.275707
(2005.0, 2015.0] 0.225000
(1.0, 1985.0]    0.150657
(0.0, 1.0]       0.000000
Name: construction_year, dtype: float64
```

```
In [78]: df_train['amount_tsh'].value_counts(normalized=True)

Out[78]: (-1, 0]          0.700993
(20, 350000]     0.249495
(0, 20]          0.049512
Name: amount_tsh, dtype: float64
```

```
In [79]: df_train['amount_tsh']

Out[79]: 0          (20, 350000]
1          (-1, 0]
2          (20, 350000]
3          (-1, 0]
4          (-1, 0]
...
59395      (0, 20]
59396      (20, 350000]
59397      (-1, 0]
59398      (-1, 0]
59399      (-1, 0]
Name: amount_tsh, Length: 59400, dtype: object
Categories (3, interval[int64]): [(-1, 0], (0, 20], (20, 350000]]
```

```
In [80]: df_train['amount_tsh'].value_counts(normalized=True)

Out[80]: (-1, 0]          0.700993
(20, 350000]     0.249495
(0, 20]          0.049512
Name: amount_tsh, dtype: float64
```

```
In [81]: df_train.columns

Out[81]: Index(['id_left', 'amount_tsh', 'longitude', 'region_code', 'district_code', 'extraction_type_group', 'manager_quality_group', 'quantity', 'source_waterpoint_type', 'installer_new_status_group'],
              dtype='object')
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

### 4.4.3 Clean Target

```
In [82]: df_train['status_group'].value_counts(nc
```

```
Out[82]: functional          0.436229
non functional          0.305522
NaN                    0.200471
functional needs repair  0.057778
Name: status_group, dtype: float64
```

```
In [83]: df_train.status_group.fillna(0, inplace=
```

```
In [84]: df_train['status_group'].value_counts(nc
```

```
Out[84]: functional          0.436229
non functional          0.305522
0                    0.200471
functional needs repair  0.057778
Name: status_group, dtype: float64
```

```
In [85]: # delete rows where prediction values dic
df_train = df_train[df_train.status_grou
```

### 4.4.4 Visualizations

```
In [184]: df_train['status_group_2'] = df_train['s
di = {'functional': 'functional', 'non f
df_train['status_group_2'].replace(di, i
```

```
In [191]: # creating a list
visuals = df_train.columns

# items to be removed
unwanted_ele = ['id_left', 'longitude', 'l

visuals = [ele for ele in visuals if ele
```

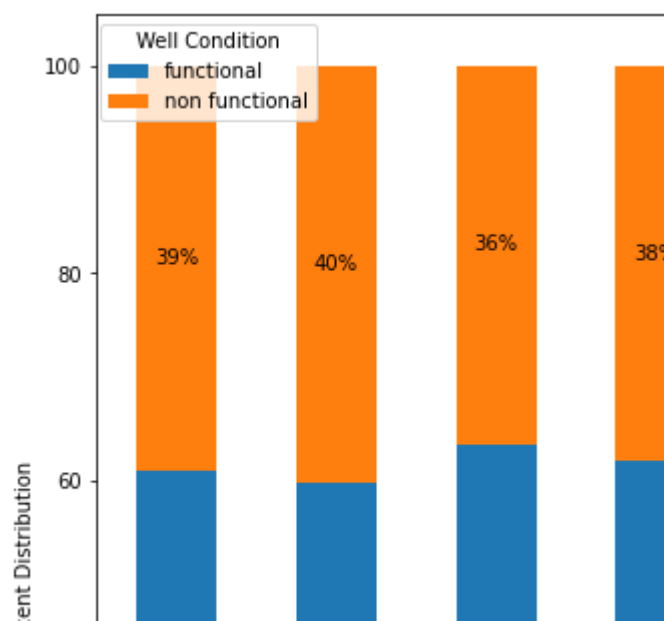
## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [198]: for column in visuals:
            ax= pd.crosstab(df_train[column], df_train['functional'],
                           ax_1 = ax.plot.bar(figsize=(10,10),
                           display(ax)
                           plt.legend(loc='upper center', bbox_
                           plt.xlabel(column)
                           plt.xticks(rotation=65)
                           plt.ylabel('Percent Distribution')

            for rec in ax_1.patches:
                height = rec.get_height()
                ax_1.text(rec.get_x() + rec.get_
                           rec.get_y() + height / 2,
                           "{:.0f}%".format(height),
                           ha='center',
                           va='bottom')

            plt.show()
```



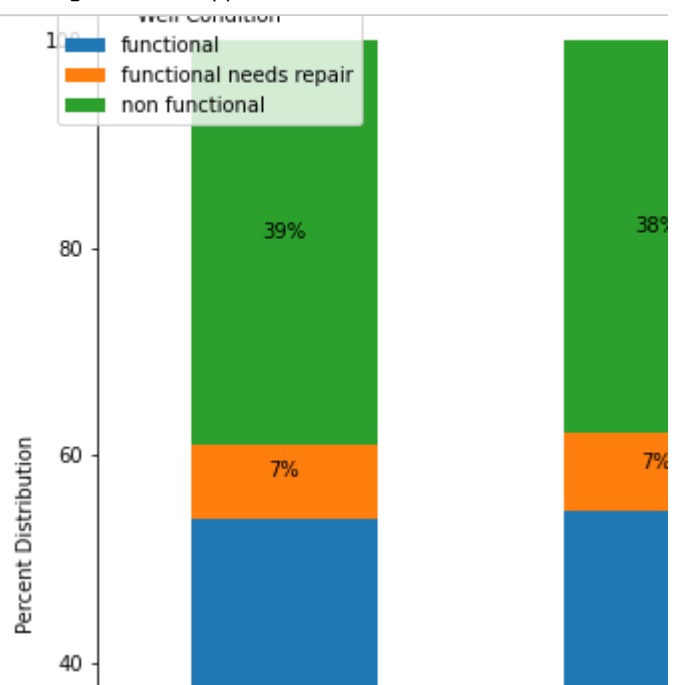
## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [199]: for column in visuals:
            ax= pd.crosstab(df_train[column], df_train['water_condition'])
            ax_1 = ax.plot.bar(figsize=(10,10),style='b')
            display(ax)
            plt.legend(loc='upper center', bbox_=(0.5, 0.95), shadow=True)
            plt.xlabel(column)
            plt.xticks(rotation=65)
            plt.ylabel('Percent Distribution')

            for rec in ax_1.patches:
                height = rec.get_height()
                ax_1.text(rec.get_x() + rec.get_width() / 2,
                        rec.get_y() + height / 2,
                        "{:.0f}%".format(height),
                        ha='center',
                        va='bottom')

            plt.show()
```



```
In [203]: df_train['latitude'].value_counts(normalized=True)
```

```
Out[203]: -2.000000e-08    0.030574
           -6.980122e+00    0.000042
           -2.494546e+00    0.000042
           -6.956746e+00    0.000042
           -6.964258e+00    0.000042
           ...
           -6.899767e+00    0.000021
           -4.422470e+00    0.000021
           -8.446126e+00    0.000021
           -3.220441e+00    0.000021
           -7.714786e+00    0.000021
           Name: latitude, Length: 45996, dtype: float64
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

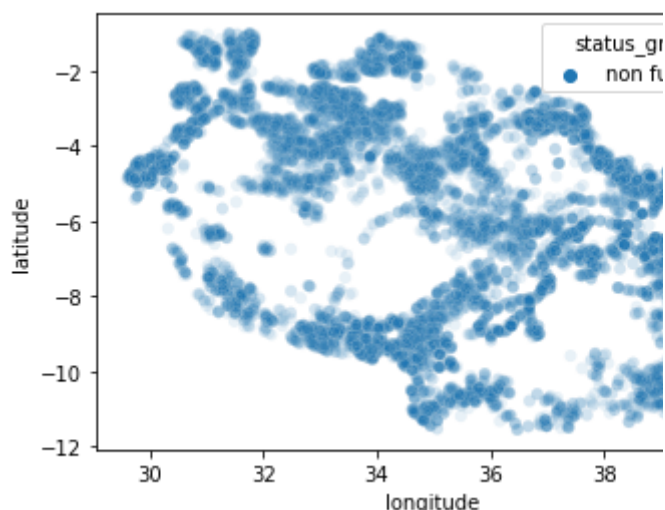
```
In [204]: df_train['longitude'].value_counts(normalized=True)
```

```
Out[204]: 0.000000    0.030574
          37.314250    0.000042
          32.984790    0.000042
          39.105307    0.000042
          39.095087    0.000042
          ...
          29.810473    0.000021
          37.383938    0.000021
          38.991929    0.000021
          38.455444    0.000021
          36.781854    0.000021
```

```
Name: longitude, Length: 45994, dtype: float64
```

```
In [205]: df_visual = df_train[df_train.longitude > 30]
```

```
In [223]: df_visual.plot(kind='scatter', x='longitude', y='latitude', c='status_group', legend=True, figsize=(15, 10))
```

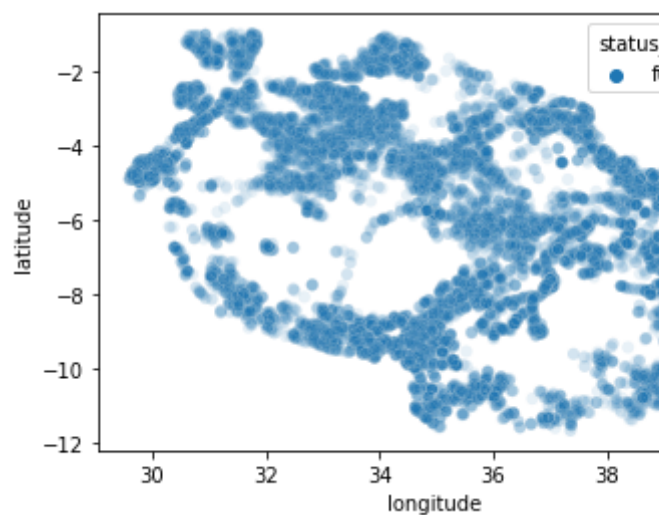




## Contents 🔄 ⚙️

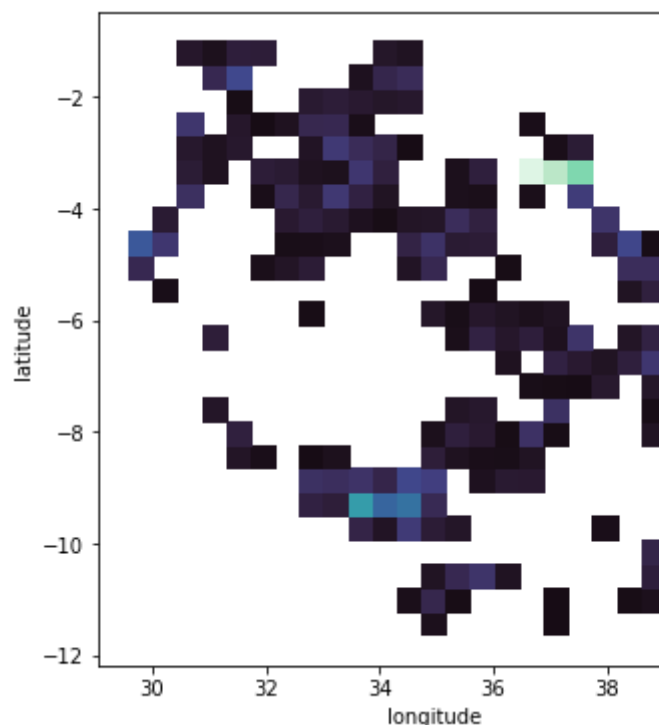
- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [224]: scatterplot('longitude', 'latitude', data=df_train)
plt.show()
```



```
In [218]: # Draw a combo histogram and scatterplot
f, ax = plt.subplots(figsize=(6, 6))
#sns.scatterplot(x='longitude', y='latitude', data=df_train)
sns.histplot(x='longitude', y='latitude',
             bins=25, pthresh=.1, cmap='magma',
             data=df_visual[df_visual['status'] == 'f'])
#sns.kdeplot(x='longitude', y='latitude', data=df_train)
```

```
Out[218]: <AxesSubplot:xlabel='longitude', ylabel='latitude'>
```



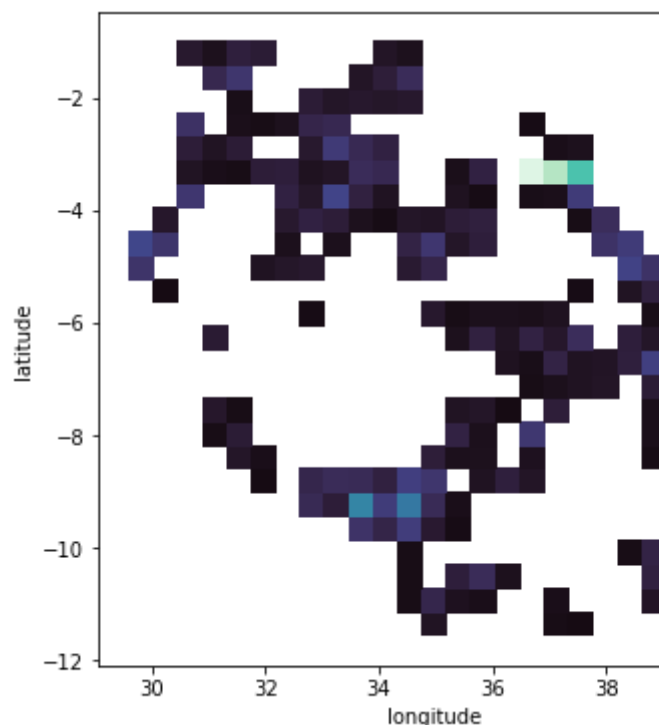


## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [219]: # Draw a combo histogram and scatterplot
f, ax = plt.subplots(figsize=(6, 6))
#sns.scatterplot(x='longitude', y='latitude')
sns.histplot(x='longitude', y='latitude',
             bins=25, pthresh=.1, cmap='magma')
#sns.kdeplot(x='longitude', y='latitude')
```

Out[219]: <AxesSubplot:xlabel='longitude', ylabel='latitude'>



## 5 Train Test Split

```
In [86]: X=df_train.drop(columns = ['status_group', 'target'])
y=df_train['status_group'] #Target
```

```
In [87]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Contents

▼ 1 Business Problem:

1.1 Background

1.2 Problem Statement

2 Import Libraries

▼ 3 Data Exploration

3.1 Data Fields

▼ 4 Data Preparation

▼ 4.1 Missing Values

4.1.1 permit

4.1.2 construction\_year

▼ 4.2 Replace misspellings and group smaller

4.2.1 Reformat Installer col

4.2.2 Reformat Funder col

▼ 4.2.3 Group Other Remaining Columns

4.2.3.1 lga

▼ 4.3 Columns to drop

4.3.1 Mostly Empty

4.3.2 Many Individual Values

▼ 4.3.3 Not Significant

4.3.3.1 The features scheme\_manage

▼ 4.4 Categorical and Numerical

4.4.1 Join Target: df\_train\_set to df\_train

4.4.2 Column Binning

4.4.3 Clean Target

4.4.4 Visualizations

5 Train Test Split

▼ 6 Encode Features

6.1 X\_train Encode

6.2 X\_test Encode

▼ 6.3 Delete 'other' columns OHE

6.3.1 Delete Other Train

6.3.2 Delete Other Test

6.3.3 Label Encode Target

▼ 7 Model Development

▼ 7.1 Random Forest Classifier

7.1.1 Check for Overfit

7.1.2 Feature Importance

7.1.3 Model reiteration - parameter tuning

▼ 7.2 Gradient Boosting Classifier

7.2.1 Check for Overfit

7.2.2 Model reiteration - parameter tuning

▼ 7.3 Logistic Regression

7.3.1 unbalanced

7.3.2 balanced

7.3.3 Check for Overfit

7.3.4 Model reiteration - parameter tuning

```
In [88]: X_train
```

id	lga	permit	construction_year	...	management_group	payment_group
53	other	1	(1985.0, 2005.0]	...	user-group	newly paved
5	hai	1	(1985.0, 2005.0]	...	user-group	other
2	other	0	(1985.0, 2005.0]	...	user-group	newly paved
3	same	1	(1985.0, 2005.0]	...	user-group	paved annually
3	other	0	(1.0, 1985.0]	...	user-group	newly paved
...	...	...	...	...	...	...
7	mbarali	1	NaN	...	user-group	newly paved

```
In [89]: X_train.shape
```

```
Out[89]: (33244, 22)
```

```
In [90]: X_test.shape
```

```
Out[90]: (14248, 22)
```

## 6 Encode Features

```
In [91]: df_train.columns
```

```
Out[91]: Index(['id_left', 'amount_tsh', 'longitude', 'region_code', 'district_code', 'extraction_type_group', 'management_group', 'quality_group', 'quantity', 'source_waterpoint_type', 'installer_new', 'status_group'], dtype='object')
```

```
In [92]: col_names = list(df_train.columns)
```

```
In [93]: col_names_ohe = col_names
```

```
In [94]: remove_ohe = ['id_left', 'id_right', 'perm
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [95]: for ele in remove_ohe:
          if ele in col_names_ohe:
              col_names_ohe.remove(ele)
          col_names_ohe
```

```
Out[95]: ['amount_tsh',
          'basin',
          'region_code',
          'district_code',
          'lga',
          'construction_year',
          'extraction_type_group',
          'management',
          'management_group',
          'payment',
          'quality_group',
          'quantity',
          'source_type',
          'source_class',
          'waterpoint_type',
          'installer_new',
          'funder_new']
```

```
In [96]: X_train[col_names] = X_train[col_names_ohe]
          X_test[col_names] = X_test[col_names_ohe]
```

```
In [97]: X_train_ohe = X_train[col_names_ohe]
          X_test_ohe = X_test[col_names_ohe]
```

## 6.1 X\_train Encode

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [98]: X\_train.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 33244 entries, 13863 to 30963
Data columns (total 22 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   id_left                33244 non-null  int64
 1   amount_tsh            33244 non-null  float64
 2   longitude              33244 non-null  float64
 3   latitude               33244 non-null  float64
 4   basin                 33244 non-null  object
 5   region_code           33244 non-null  object
 6   district_code         33244 non-null  object
 7   lga                   33244 non-null  object
 8   permit                33244 non-null  object
 9   construction_year     21674 non-null  int64
10   extraction_type_group 33244 non-null  object
11   management            33244 non-null  object
12   management_group      33244 non-null  object
13   payment               33244 non-null  object
14   quality_group         33244 non-null  object
15   quantity              33244 non-null  int64
16   source_type           33244 non-null  object
17   source_class          33244 non-null  object
18   waterpoint_type       33244 non-null  object
19   installer_new         33244 non-null  object
20   funder_new            33244 non-null  object
21   id_right              33244 non-null  int64
dtypes: category(17), float64(3), int64(2)
memory usage: 2.1 MB
```

Change Sci kit learn OHE set sparse to false set handle\_unknown='ignore'

In [99]: enc = OneHotEncoder(handle\_unknown='ignore')

In [100]: enc.fit(X\_train\_ohe)

Out[100]: OneHotEncoder(handle\_unknown='ignore', sparse=False)

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [101]: enc.categories_
```

```
Out[101]: [array([Interval(-1, 0, closed='right'),
              Interval(20, 350000, closed='right')],
            array(['Internal', 'Lake Nyasa', 'Lake
              Lake Victoria', 'Pangani', 'Rufiji',
              'Wami / Ruvu'], dtype=object),
            array([ 1,  2,  3,  4,  5,  6,  7,  8,
                  18, 19, 20, 21, 24, 40, 60, 80],
            array([ 0,  1,  2,  3,  4,  5,  6,  7,
                  63, 67, 80]),
            array(['arusha rural', 'bagamoyo', 'bahi
              kahama', 'karagwe', 'kasulu', 'kibaha',
              'kilombero', 'kilosa', 'kwimba', 'makete',
              'maswa', 'mbarali', 'mbozi', 'moshi rural',
              'mpanda', 'mvomeri', 'other', 'rombo',
              'rungwe', 'sangha', 'songea rural', 'ulanga'],
              dtype=object),
            array([Interval(1.0, 1985.0, closed='right'),
                  Interval(1985.0, 2005.0, closed='right'),
                  Interval(2005.0, 2015.0, closed='right')],
              dtype=object)]
```

```
In [102]: X_train2 = enc.transform(X_train_ohe)
```

```
In [103]: X_train2.shape
```

```
Out[103]: (33244, 197)
```

```
In [104]: enc.get_feature_names()
column_name = enc.get_feature_names()
one_hot_encoded_frame = pd.DataFrame(X_train2, columns=column_name)
```

```
In [105]: X_train = one_hot_encoded_frame.reset_index(drop=True)
```

```
In [106]: X_train.drop(col_names_ohe, axis=1, inplace=True)
```

Out[106]:

	x0_(-1, 0]	x0_(0, 20]	x0_(20, 350000]	x1_Internal	x1_Lake Nyasa
0	1.0	0.0	0.0	0.0	0.0
1	0.0	0.0	1.0	0.0	0.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	1.0	0.0	0.0
4	1.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...
33239	1.0	0.0	0.0	0.0	0.0
33240	0.0	1.0	0.0	0.0	0.0
33241	1.0	0.0	0.0	0.0	0.0
33242	0.0	0.0	1.0	0.0	0.0

```
In [107]: del X_train['id_left']
           del X_train['id_right']
```

```
In [108]: X_train
```

Out[108]:

	x0_(-1, 0]	x0_(0, 20]	x0_(20, 350000]	x1_Internal	x1_Lake Nyasa
0	1.0	0.0	0.0	0.0	0.0
1	0.0	0.0	1.0	0.0	0.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	1.0	0.0	0.0
4	1.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...
33239	1.0	0.0	0.0	0.0	0.0
33240	0.0	1.0	0.0	0.0	0.0
33241	1.0	0.0	0.0	0.0	0.0
33242	0.0	0.0	1.0	0.0	0.0

## 6.2 X\_test Encode

### Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [109]: X\_test.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 14248 entries, 50633 to 3017
Data columns (total 22 columns):
#   Column                                Non-Null Count
---  ---
0   id_left                               14248 non-null
1   amount_tsh                           14248 non-null
2   longitude                             14248 non-null
3   latitude                              14248 non-null
4   basin                                 14248 non-null
5   region_code                           14248 non-null
6   district_code                         14248 non-null
7   lga                                   14248 non-null
8   permit                               14248 non-null
9   construction_year                     9298 non-null
10  extraction_type_group                 14248 non-null
11  management                           14248 non-null
12  management_group                     14248 non-null
13  payment                              14248 non-null
..   ..
```

In [110]: enc.categories\_

```
Out[110]: [array([Interval(-1, 0, closed='right'),
          Interval(20, 350000, closed='right')],
          array(['Internal', 'Lake Nyasa', 'Lake
          'Lake Victoria', 'Pangani', 'Rufiji',
          'Wami / Ruvu'], dtype=object),
          array([ 1,  2,  3,  4,  5,  6,  7,  8,
                  18, 19, 20, 21, 24, 40, 60, 80],
          array([ 0,  1,  2,  3,  4,  5,  6,  7,
                  63, 67, 80]),
          array(['arusha rural', 'bagamoyo', 'barotse',
          'kahama', 'karagwe', 'kasulu', 'kibaha',
          'kilombero', 'kilosa', 'kwimba', 'makete',
          'maswa', 'mbarali', 'mbarali', 'moshi rural',
          'mpanda', 'mvomeri', 'other', 'rombo', 'rungwe',
          'sarungu', 'songea rural', 'ulanga'], dtype=object),
          array([Interval(1.0, 1985.0, closed='right'),
          Interval(1985.0, 2005.0, closed='right'),
          Interval(2005.0, 2015.0, closed='right')])]
```

In [111]: X\_test2 = enc.transform(X\_test\_ohe)

In [112]: X\_test2.shape

Out[112]: (14248, 197)

```
In [113]: enc.get_feature_names()
column_name = enc.get_feature_names()
one_hot_encoded_frame = pd.DataFrame(X_test2,
```



## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [114]: X_test = one_hot_encoded_frame.reset_index(drop=True)
```

```
In [115]: X_test.drop(col_names_ohe, axis=1, inplace=True)
X_test
```

Out[115]:

	x0_(-1, 0]	x0_(0, 20]	x0_(20, 350000]	x1_Internal	x1_Lake Nyasa
0	1.0	0.0	0.0	0.0	0.0
1	1.0	0.0	0.0	0.0	1.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	1.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0
...	...	...	...	...	...
14243	1.0	0.0	0.0	0.0	0.0
14244	1.0	0.0	0.0	1.0	0.0
14245	1.0	0.0	0.0	1.0	0.0
14246	0.0	0.0	1.0	0.0	0.0

```
In [116]: del X_test['id_left']
del X_test['id_right']
```

```
In [117]: X_test
```

Out[117]:

	x0_(-1, 0]	x0_(0, 20]	x0_(20, 350000]	x1_Internal	x1_Lake Nyasa
0	1.0	0.0	0.0	0.0	0.0
1	1.0	0.0	0.0	0.0	1.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	1.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0
...	...	...	...	...	...
14243	1.0	0.0	0.0	0.0	0.0
14244	1.0	0.0	0.0	1.0	0.0
14245	1.0	0.0	0.0	1.0	0.0
14246	0.0	0.0	1.0	0.0	0.0

```
In [118]: X_train.shape
```

Out[118]: (33244, 200)



```
In [119]: X_test.shape
```

```
Out[119]: (14248, 200)
```

## 6.3 Delete 'other' columns OHE

For Text data which was too low incidence to make it

### 6.3.1 Delete Other Train

```
In [120]: X_train
```

```
Out[120]:
```

	x0_(-1, 0]	x0_(0, 20]	x0_(20, 350000]	x1_Internal	x1_Lake Nyasa
0	1.0	0.0	0.0	0.0	0.0
1	0.0	0.0	1.0	0.0	0.0
2	1.0	0.0	0.0	0.0	0.0
3	0.0	0.0	1.0	0.0	0.0
4	1.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...
33239	1.0	0.0	0.0	0.0	0.0
33240	0.0	1.0	0.0	0.0	0.0
33241	1.0	0.0	0.0	0.0	0.0
33242	0.0	0.0	1.0	0.0	0.0
33243	1.0	0.0	0.0	1.0	0.0

33244 rows × 200 columns

```
In [121]: other_cols = [col for col in X_train.colnames() if col not in ['x0', 'x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10', 'x11', 'x12', 'x13', 'x14', 'x15', 'x16', 'x17', 'x18', 'x19', 'x20', 'x21', 'x22', 'x23', 'x24', 'x25', 'x26', 'x27', 'x28', 'x29', 'x30', 'x31', 'x32', 'x33', 'x34', 'x35', 'x36', 'x37', 'x38', 'x39', 'x40', 'x41', 'x42', 'x43', 'x44', 'x45', 'x46', 'x47', 'x48', 'x49', 'x50', 'x51', 'x52', 'x53', 'x54', 'x55', 'x56', 'x57', 'x58', 'x59', 'x60', 'x61', 'x62', 'x63', 'x64', 'x65', 'x66', 'x67', 'x68', 'x69', 'x70', 'x71', 'x72', 'x73', 'x74', 'x75', 'x76', 'x77', 'x78', 'x79', 'x80', 'x81', 'x82', 'x83', 'x84', 'x85', 'x86', 'x87', 'x88', 'x89', 'x90', 'x91', 'x92', 'x93', 'x94', 'x95', 'x96', 'x97', 'x98', 'x99']]
```

```
Out[121]: ['x4_other',
            'x6_other',
            'x6_other handpump',
            'x6_other motorpump',
            'x7_other',
            'x7_other - school',
            'x8_other',
            'x9_other',
            'x12_other',
            'x14_other',
            'x15_other',
            'x16_other']
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [122]: X_train = X_train.drop(other_cols, axis=
X_train.columns
```

```
Out[122]: Index(['x0_(-1, 0]', 'x0_(0, 20]', 'x0_
'x1_Lake Nyasa', 'x1_Lake Rukwa',
'x1_Lake Victoria', 'x1_Pangani',
...
'x16_private individual', 'x16_rv
'x16_unicef', 'x16_world bank', '
'latitude', 'permit'],
dtype='object', length=188)
```

### 6.3.2 Delete Other Test

```
In [123]: X_test = X_test.drop(other_cols, axis=1)
X_test.columns
```

```
Out[123]: Index(['x0_(-1, 0]', 'x0_(0, 20]', 'x0_
'x1_Lake Nyasa', 'x1_Lake Rukwa',
'x1_Lake Victoria', 'x1_Pangani',
...
'x16_private individual', 'x16_rv
'x16_unicef', 'x16_world bank', '
'latitude', 'permit'],
dtype='object', length=188)
```

```
In [124]: other_cols_test = [col for col in X_test
other_cols_test
```

```
Out[124]: []
```

### 6.3.3 Label Encode Target

```
In [125]: di = {0: 0, 'functional': 1, 'non functi
y_train.replace(di, inplace=True)
```

```
In [126]: di = {0: 0, 'functional': 1, 'non functi
y_test.replace(di, inplace=True)
```

## 7 Model Development

### 7.1 Random Forest Classifier

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [127]: # Instantiate and fit the RandomForestClassifier
rfc=RandomForestClassifier(n_estimators=100)
rfc.fit(X_train,y_train)
```

```
Out[127]: RandomForestClassifier(class_weight='balanced')
```

```
In [128]: training_start = time.perf_counter()
rfc.fit(X_train, y_train)
training_end = time.perf_counter()
prediction_start = time.perf_counter()
y_pred = rfc.predict(X_test)
y_pred_train = rfc.predict(X_train)
prediction_end = time.perf_counter()
acc_rfc = (y_pred == y_test).sum().astype(float)
rfc_train_time = training_end-training_start
rfc_prediction_time = prediction_end-prediction_start
print("Scikit-Learn's Random Forest Classifier")
print("Time consumed for training: %4.3f seconds" % rfc_train_time)
print("Time consumed for prediction: %6.3f seconds" % rfc_prediction_time)
```

```
Scikit-Learn's Random Forest Classifier'
Time consumed for training: 1.057 seconds
Time consumed for prediction: 0.24293 seconds
```

```
In [129]: print(metrics.classification_report(y_test, y_pred))
```

	precision	recall	f1-score
1	0.55	0.65	0.60
2	0.39	0.32	0.35
3	0.10	0.06	0.08
accuracy			0.50
macro avg	0.35	0.34	0.34
weighted avg	0.46	0.48	0.47

### 7.1.1 Check for Overfit

```
In [130]: # View confusion matrix for train data and test data
confusion_matrix(y_train, y_pred_train)
```

```
Out[130]: array([[17717, 294, 113],
 [ 865, 11780, 70],
 [ 207, 116, 2082]])
```

```
In [131]: # View confusion matrix for test data and train data
confusion_matrix(y_test, y_pred)
```

```
Out[131]: array([[5026, 2429, 333],
 [3449, 1760, 224],
 [ 638, 329, 60]])
```

Contents

▼ 1 Business Problem:

1.1 Background

1.2 Problem Statement

2 Import Libraries

▼ 3 Data Exploration

3.1 Data Fields

▼ 4 Data Preparation

▼ 4.1 Missing Values

4.1.1 permit

4.1.2 construction\_year

▼ 4.2 Replace misspellings and group smaller

4.2.1 Reformat Installer col

4.2.2 Reformat Funder col

▼ 4.2.3 Group Other Remaining Columns

4.2.3.1 lga

▼ 4.3 Columns to drop

4.3.1 Mostly Empty

4.3.2 Many Individual Values

▼ 4.3.3 Not Significant

4.3.3.1 The features scheme\_manage

▼ 4.4 Categorical and Numerical

4.4.1 Join Target: df\_train\_set to df\_train

4.4.2 Column Binning

4.4.3 Clean Target

4.4.4 Visualizations

5 Train Test Split

▼ 6 Encode Features

6.1 X\_train Encode

6.2 X\_test Encode

▼ 6.3 Delete 'other' columns OHE

6.3.1 Delete Other Train

6.3.2 Delete Other Test

6.3.3 Label Encode Target

▼ 7 Model Development

▼ 7.1 Random Forest Classifier

7.1.1 Check for Overfit

7.1.2 Feature Importance

7.1.3 Model reiteration - parameter tuning

▼ 7.2 Gradient Boosting Classifier

7.2.1 Check for Overfit

7.2.2 Model reiteration - parameter tuning

▼ 7.3 Logistic Regression

7.3.1 unbalanced

7.3.2 balanced

7.3.3 Check for Overfit

7.3.4 Model reiteration - parameter tuning

```
In [132]: rfc_training_preds = rfc.predict(X_train)
rfc_training_accuracy = accuracy_score(y_train, rfc_training_preds)

rfc_val_preds = rfc.predict(X_test) # y_val
rfc_val_accuracy = accuracy_score(y_test, rfc_val_preds)
print(rfc_training_accuracy)
print(rfc_val_accuracy)

0.9499157742750571
0.48048848961257723
```

Because we are looking for the minority class accuracy

Lets take a look at other metrics in the classification report

```
In [133]: print(metrics.classification_report(y_train, rfc_val_preds))
```

	precision	recall	f1-score
1	0.94	0.98	(
2	0.97	0.93	(
3	0.92	0.87	(
accuracy			(
macro avg	0.94	0.92	(
weighted avg	0.95	0.95	(

```
In [134]: print(metrics.classification_report(y_test, rfc_val_preds))
```

	precision	recall	f1-score
1	0.55	0.65	(
2	0.39	0.32	(
3	0.10	0.06	(
accuracy			(
macro avg	0.35	0.34	(
weighted avg	0.46	0.48	(

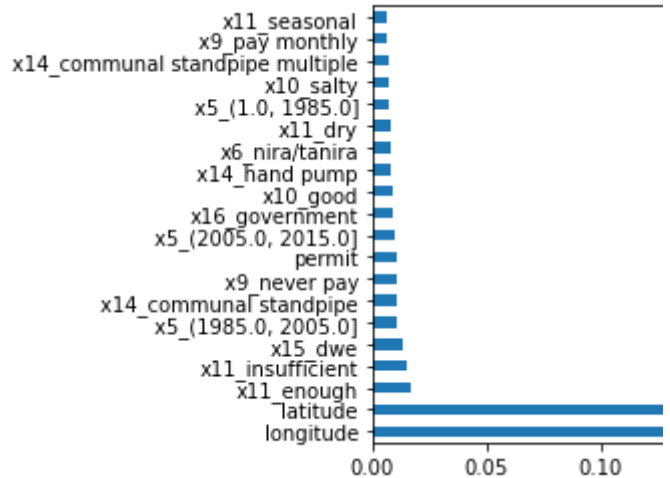
Add labels so it's clear which class is being described

### 7.1.2 Feature Importance

```
In [135]: # create list so that random forest model can use it
rfc_columns = list(X_train.columns)
```

```
In [136]: feat_importances = pd.Series(rfc.feature
feat_importances.nlargest(20).plot(kind=
```

```
Out[136]: <AxesSubplot:>
```



## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

### 7.1.3 Model reiteration - parameter tuning

Don't set anything to random, choose the ones that you want to test a percentage of features.

```
In [137]: model_params = {
# number of trees
'n_estimators': [5,10,20,50,100,200]
# number of max features
'max_features': [10,15,20,50,100],
# max number of levels in each decision tree
'max_depth': [5,10,15],
'min_samples_split' : [100,1000]
}
```

```
In [138]: # Instantiate and fit the RandomForestClassifier
rfc=RandomForestClassifier(class_weight='balanced')
rfc = RandomizedSearchCV(rfc,model_params,cv=3)
rfc.fit(X_train,y_train)
```

```
Out[138]: RandomizedSearchCV(cv=3,
estimator=RandomForestClassifier(class_weight='balanced'),
param_distributions={'max_depth': [5, 10, 15],
'n_estimators': [5, 10, 20, 50, 100, 200],
'max_features': [10, 15, 20, 50, 100],
'min_samples_split': [100, 1000]},
random_state=1)
```

## Contents 🔍 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [139]: print(rfc.best_estimator_.get_params())
```

```
{'bootstrap': True, 'ccp_alpha': 0.0, 'criterion': 'gini', 'max_depth': 50, 'max_leaf_nodes': None, 'max_samples': 1.0, 'min_samples_leaf': 1, 'min_samples_split': 10, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 100, 'oob_score': False, 'random_state': None, 'verbose': 0}
```

```
In [140]: print(metrics.classification_report(y_train, y_pred_train))
```

	precision	recall	f1-score
1	0.94	0.98	0.96
2	0.97	0.93	0.95
3	0.92	0.87	0.89
accuracy			0.95
macro avg	0.94	0.92	0.93
weighted avg	0.95	0.95	0.95

```
In [141]: print(metrics.classification_report(y_test, y_pred_test))
```

	precision	recall	f1-score
1	0.55	0.65	0.60
2	0.39	0.32	0.35
3	0.10	0.06	0.08
accuracy			0.34
macro avg	0.35	0.34	0.34
weighted avg	0.46	0.48	0.47

## 7.2 Gradient Boosting Classifier

```
In [142]: gb_clf = GradientBoostingClassifier()
gb_clf.fit(X_train, y_train)
y_pred = gb_clf.predict(X_test)
y_pred_train = gb_clf.predict(X_train)

# print("Accuracy score (training): {0:.4f}".format(metrics.accuracy_score(y_train, y_pred_train)))
# print("Accuracy score (validation): {0:.4f}".format(metrics.accuracy_score(y_test, y_pred)))
```

Contents

▼ 1 Business Problem:

1.1 Background

1.2 Problem Statement

2 Import Libraries

▼ 3 Data Exploration

3.1 Data Fields

▼ 4 Data Preparation

▼ 4.1 Missing Values

4.1.1 permit

4.1.2 construction\_year

▼ 4.2 Replace misspellings and group smaller

4.2.1 Reformat Installer col

4.2.2 Reformat Funder col

▼ 4.2.3 Group Other Remaining Columns

4.2.3.1 lga

▼ 4.3 Columns to drop

4.3.1 Mostly Empty

4.3.2 Many Individual Values

▼ 4.3.3 Not Significant

4.3.3.1 The features scheme\_manage

▼ 4.4 Categorical and Numerical

4.4.1 Join Target: df\_train\_set to df\_train

4.4.2 Column Binning

4.4.3 Clean Target

4.4.4 Visualizations

5 Train Test Split

▼ 6 Encode Features

6.1 X\_train Encode

6.2 X\_test Encode

▼ 6.3 Delete 'other' columns OHE

6.3.1 Delete Other Train

6.3.2 Delete Other Test

6.3.3 Label Encode Target

▼ 7 Model Development

▼ 7.1 Random Forest Classifier

7.1.1 Check for Overfit

7.1.2 Feature Importance

7.1.3 Model reiteration - parameter tuning

▼ 7.2 Gradient Boosting Classifier

7.2.1 Check for Overfit

7.2.2 Model reiteration - parameter tuning

▼ 7.3 Logistic Regression

7.3.1 unbalanced

7.3.2 balanced

7.3.3 Check for Overfit

7.3.4 Model reiteration - parameter tuning

```
In [143]: print(metrics.classification_report(y_train, y_pred_train))
```

	precision	recall	f1-score
1	0.55	0.99	0.72
2	0.34	0.01	0.06
3	0.50	0.00	0.00
accuracy			0.50
macro avg	0.46	0.33	0.43
weighted avg	0.46	0.54	0.50

This model did worse at finding the minority class

### 7.2.1 Check for Overfit

```
In [144]: # View confusion matrix for train data and train set
confusion_matrix(y_train, y_pred_train)
```

```
Out[144]: array([[18074,    50,     0],
                 [12410,   305,     0],
                 [ 2360,    32,    13]])
```

```
In [145]: # View confusion matrix for test data and test set
confusion_matrix(y_test, y_pred)
```

```
Out[145]: array([[7706,    82,     0],
                 [5386,    46,     1],
                 [1019,     7,     1]])
```

```
In [146]: print(metrics.classification_report(y_train, y_pred_train))
```

	precision	recall	f1-score
1	0.55	1.00	0.72
2	0.79	0.02	0.11
3	1.00	0.01	0.02
accuracy			0.50
macro avg	0.78	0.34	0.43
weighted avg	0.67	0.55	0.50



## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [147]: `print(metrics.classification_report(y_test, y_pred))`

	precision	recall	f1-score
1	0.55	0.99	(0.72)
2	0.34	0.01	(0.06)
3	0.50	0.00	(0.00)
accuracy			(0.74)
macro avg	0.46	0.33	(0.40)
weighted avg	0.46	0.54	(0.50)

## 7.2.2 Model reiteration - parameter tuning

In [148]: `model_params = {`  
           `# number of boosting stages`  
           `'n_estimators': [5,10,20,50,100,200]`  
           `# number of max features`  
           `'max_features': [10,15,20,50,100],`  
           `# Learning rate`  
           `'learning_rate': [.25,.5,.75,1],`  
           `#The minimum number of samples required to split an internal node`  
           `'min_samples_split' : [100,1000]`  
           `}`

In [149]: `gb_clf = RandomizedSearchCV(gb_clf,model_params,cv=3,random_state=1)`  
           `gb_clf.fit(X_train,y_train)`

Out[149]: `RandomizedSearchCV(cv=3, estimator=GradientBoostingClassifier(),`  
           `param_distributions={'learning_rate': [0.25, 0.5, 0.75, 1], 'max_features': [10, 15, 20, 50, 100], 'min_samples_split': [100, 1000], 'n_estimators': [5, 10, 20, 50, 100, 200]},`

`random_state=1)`

In [150]: `print(metrics.classification_report(y_test, y_pred))`

	precision	recall	f1-score
1	0.55	1.00	(0.72)
2	0.79	0.02	(0.12)
3	1.00	0.01	(0.04)
accuracy			(0.74)
macro avg	0.78	0.34	(0.56)
weighted avg	0.67	0.55	(0.60)

```
In [151]: print(metrics.classification_report(y_te
```

	precision	recall	f1-score
1	0.55	0.99	(
2	0.34	0.01	(
3	0.50	0.00	(
accuracy			(
macro avg	0.46	0.33	(
weighted avg	0.46	0.54	(

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [152]: lr_list = [0.05, 0.075, 0.1, 0.25, 0.5,
                    1.0]

for learning_rate in lr_list:
    gb_clf = GradientBoostingClassifier()
    gb_clf.fit(X_train, y_train)

    print("Learning rate: ", learning_rate)
    print("Accuracy score (training): {}".format(
        accuracy_score(X_train, gb_clf.predict(X_train))
    ))
    print("Accuracy score (validation): {}".format(
        accuracy_score(X_test, gb_clf.predict(X_test))
    ))

    print("Accuracy score (training): 0.546")
    print("Accuracy score (validation): 0.546")
    print("Learning rate: 0.075")
    print("Accuracy score (training): 0.546")
    print("Accuracy score (validation): 0.546")
    print("Learning rate: 0.1")
    print("Accuracy score (training): 0.547")
    print("Accuracy score (validation): 0.546")
    print("Learning rate: 0.25")
    print("Accuracy score (training): 0.550")
    print("Accuracy score (validation): 0.545")
    print("Learning rate: 0.5")
    print("Accuracy score (training): 0.558")
    print("Accuracy score (validation): 0.541")
    print("Learning rate: 0.75")
    print("Accuracy score (training): 0.558")
    print("Accuracy score (validation): 0.539")
    print("Learning rate: 1")
    print("Accuracy score (training): 0.561")
    print("Accuracy score (validation): 0.538")
```

## 7.3 Logistic Regression

### 7.3.1 unbalanced

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [153]: # Logistic model
log_clf = LogisticRegression(random_state=0)
log_model = log_clf.fit(X_train, y_train)
log_training_preds = log_clf.predict(X_train)
log_training_accuracy = accuracy_score(y_train, log_training_preds)
log_val_preds = log_clf.predict(X_test)
log_val_accuracy = accuracy_score(y_test, log_val_preds)
```

```
In [154]: #Confusion matrix for Logistic Regression
log_matrix = confusion_matrix(y_test, log_val_preds)
print('Confusion Matrix:\n', log_matrix)
```

```
Confusion Matrix:
[[7677  111    0]
 [5361   72    0]
 [1009   18    0]]
```

```
In [155]: print(log_training_accuracy)
print(log_val_accuracy)
```

```
0.5455119720851883
0.5438658057271196
```

```
In [156]: y_pred = log_model.predict(X_test)
```

```
In [157]: print(metrics.classification_report(y_test, y_pred))
```

	precision	recall	f1-score
1	0.55	0.99	(0.72)
2	0.36	0.01	(0.07)
3	0.00	0.00	(0.00)
accuracy			(0.54)
macro avg	0.30	0.33	(0.31)
weighted avg	0.44	0.54	(0.49)

### 7.3.2 balanced

```
In [158]: # Logistic model
log_clf = LogisticRegression(multi_class='multinomial', solver='lbfgs')
log_model = log_clf.fit(X_train, y_train)
y_pred_train = log_model.predict(X_train)
y_pred = log_model.predict(X_test)
log_training_accuracy = accuracy_score(y_train, y_pred_train)
log_val_accuracy = accuracy_score(y_test, y_pred)
```

## Contents 🔄 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [159]: #Confusion matrix for Logistic Regression
log_matrix = confusion_matrix(y_test, y_pred)
print('Confusion Matrix:\n', log_matrix)
```

```
Confusion Matrix:
[[2546 2399 2843]
 [1735 1657 2041]
 [ 347  297  383]]
```

```
In [160]: print(log_training_accuracy)
print(log_val_accuracy)
```

```
0.3494465166646613
0.32186973610331276
```

```
In [161]: print(metrics.classification_report(y_test, y_pred))
```

	precision	recall	f1-score
1	0.55	0.33	0.42
2	0.38	0.30	0.34
3	0.07	0.37	0.12
accuracy			0.33
macro avg	0.33	0.33	0.33
weighted avg	0.45	0.32	0.38

### 7.3.3 Check for Overfit

```
In [162]: # View confusion matrix for train data
confusion_matrix(y_train, y_pred_train)
```

```
Out[162]: array([[6334, 5293, 6497],
                [3973, 4206, 4536],
                [ 686,  642, 1077]])
```

```
In [163]: # View confusion matrix for test data
confusion_matrix(y_test, y_pred)
```

```
Out[163]: array([[2546, 2399, 2843],
                [1735, 1657, 2041],
                [ 347,  297,  383]])
```

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- 3 Data Exploration
  - 3.1 Data Fields
- 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group smaller
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

In [164]: `print(metrics.classification_report(y_train, y_pred_train))`

	precision	recall	f1-score
1	0.58	0.35	0.44
2	0.41	0.33	0.37
3	0.09	0.45	0.16
accuracy			0.38
macro avg	0.36	0.38	0.37
weighted avg	0.48	0.35	0.41

In [165]: `print(metrics.classification_report(y_test, y_pred_test))`

	precision	recall	f1-score
1	0.55	0.33	0.42
2	0.38	0.30	0.34
3	0.07	0.37	0.12
accuracy			0.33
macro avg	0.33	0.33	0.33
weighted avg	0.45	0.32	0.38

### 7.3.4 Model reiteration - parameter tuning

In [166]: `# Logistic model`  
`log_clf = LogisticRegression(multi_class='multinomial', solver='lbfgs')`  
`log_model = log_clf.fit(X_train, y_train)`  
`y_pred_train = log_model.predict(X_train)`  
`y_pred = log_model.predict(X_test)`  
`log_training_accuracy = accuracy_score(y_train, y_pred_train)`  
`log_val_accuracy = accuracy_score(y_test, y_pred)`

In [167]: `model_params = {`  
 `#Algorithm to use in the optimization`  
 `'solver': ['newton-cg', 'sag', 'saga'],`  
`}`

In [168]: `# Instantiate and fit the LogisticRegression model`  
`log_clf = LogisticRegression(multi_class='multinomial', solver='lbfgs')`  
`log_clf = RandomizedSearchCV(log_clf, model_params, cv=5, n_iter=10)`  
`log_clf.fit(X_train, y_train)`

Out[168]: `RandomizedSearchCV(cv=3,`  
 `estimator=LogisticRegression(multi_class='multinomial', solver='lbfgs'),`  
 `param_distributions={'solver': ['newton-cg', 'sag', 'saga']},`

## Contents

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
    - ▼ 4.2.3 Group Other Remaining Columns
      - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
    - ▼ 4.3.3 Not Significant
      - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

```
In [169]: print(log_clf.best_estimator_.get_params(
{'C': 1.0, 'class_weight': 'balanced', '
'max_iter': 100, 'multi_class': 'multin
'tol': 0.0001, 'verbose': 0, 'warm_start
```

```
In [170]: print(metrics.classification_report(y_train,
```

	precision	recall	f1-score
1	0.57	0.30	(
2	0.41	0.35	(
3	0.09	0.46	(
accuracy			(
macro avg	0.36	0.37	(
weighted avg	0.47	0.33	(

```
In [171]: print(metrics.classification_report(y_test,
```

	precision	recall	f1-score
1	0.54	0.28	(
2	0.38	0.33	(
3	0.07	0.38	(
accuracy			(
macro avg	0.33	0.33	(
weighted avg	0.44	0.30	(

## 8 Conclusions

Most Important Factors in water pump functionality:

- geographic location
  - Water Quantity
  - Construction Year
  - Whether it was permitted
  - Who funded the pump
  - Who installed the pump
- 
- Installing pumps in low water quantity areas lead
  - Permitted pumps were more likely to be functional
  - Pumps older than 1985 have poor functionality

## 9 Business Recommendation

- Expand permitting system. Waterpoints which we
- Areas which already have low water quantities ne
- Add a data feature for best guess of when pump
  - This can be used to better predict lifetime of
- Review all pumps older than 1985 as these are n

## Contents 🔍 ⚙️

- ▼ 1 Business Problem:
  - 1.1 Background
  - 1.2 Problem Statement
- 2 Import Libraries
- ▼ 3 Data Exploration
  - 3.1 Data Fields
- ▼ 4 Data Preparation
  - ▼ 4.1 Missing Values
    - 4.1.1 permit
    - 4.1.2 construction\_year
  - ▼ 4.2 Replace misspellings and group small
    - 4.2.1 Reformat Installer col
    - 4.2.2 Reformat Funder col
  - ▼ 4.2.3 Group Other Remaining Columns
    - 4.2.3.1 lga
  - ▼ 4.3 Columns to drop
    - 4.3.1 Mostly Empty
    - 4.3.2 Many Individual Values
  - ▼ 4.3.3 Not Significant
    - 4.3.3.1 The features scheme\_manage
  - ▼ 4.4 Categorical and Numerical
    - 4.4.1 Join Target: df\_train\_set to df\_train
    - 4.4.2 Column Binning
    - 4.4.3 Clean Target
    - 4.4.4 Visualizations
- 5 Train Test Split
- ▼ 6 Encode Features
  - 6.1 X\_train Encode
  - 6.2 X\_test Encode
  - ▼ 6.3 Delete 'other' columns OHE
    - 6.3.1 Delete Other Train
    - 6.3.2 Delete Other Test
    - 6.3.3 Label Encode Target
- ▼ 7 Model Development
  - ▼ 7.1 Random Forest Classifier
    - 7.1.1 Check for Overfit
    - 7.1.2 Feature Importance
    - 7.1.3 Model reiteration - parameter tuning
  - ▼ 7.2 Gradient Boosting Classifier
    - 7.2.1 Check for Overfit
    - 7.2.2 Model reiteration - parameter tuning
  - ▼ 7.3 Logistic Regression
    - 7.3.1 unbalanced
    - 7.3.2 balanced
    - 7.3.3 Check for Overfit
    - 7.3.4 Model reiteration - parameter tuning

## 10 Future Work

Things to work on if more time:

- Exploration of other numerical features such as c
  - when we examined a boxplot of these features
  - However we may have removed these features
- Look more closely at pump functionality by geog
  - The latitude and longitude of the waterpoint
  - Since this is by far our best feature further in
-