# phase-3-project-jupyter-notebook

May 23, 2024

## 0.1 Final Project Submission

Please fill out: * Student name: Alex Kipkurui korir * Student pace: part time * Scheduled project review date/time: 13th May 2024 to 22nd May 2024 * Instructor name: William Okomba/Noah Kandie * Blog post URL:

## 0.2 Business Understanding

### 0.2.1 Introduction:

Tanzania, a country with a population exceeding 57 million, faces significant challenges in providing clean water to its residents. While numerous water wells have been established, many of these wells are either in disrepair or have failed entirely. Ensuring the functionality of these water wells is crucial for public health, agriculture, and overall quality of life. Predictive analytics can play a pivotal role in identifying which wells are likely to fail, need repair, or are functioning well, thus enabling proactive maintenance and efficient resource allocation

### 0.2.2 Stakeholders and Usage:

**Government of Tanzania:**

1. Objective: Improve water supply infrastructure and resource planning.
2. Usage: By analyzing patterns in well failures, the government can develop more effective strategies for constructing new wells, maintaining existing ones, and optimizing resource allocation. This can lead to better-informed decisions on where to invest in infrastructure improvements and preventative maintenance.

**Non-Governmental Organizations (NGOs):**

1. Objective: Enhance the efficiency and impact of water-related aid programs.
2. Usage: NGOs can use predictive models to prioritize wells that need urgent repairs or are at risk of failing. This enables them to deploy their resources more effectively, ensuring that their interventions have the maximum positive impact on communities reliant on these water sources.

**Local Communities:**

1. Objective: Gain reliable access to clean water.
2. Usage: By participating in data collection and reporting well conditions, local communities can contribute to the ongoing monitoring and maintenance efforts. This collaboration can

help ensure that issues are addressed promptly, minimizing the time residents are without clean water.

### 0.2.3 Conclusion:

Developing a classifier to predict the condition of water wells in Tanzania holds significant potential for improving water supply reliability across the country. By leveraging data analytics, stakeholders such as the government, NGOs, and local communities can make informed decisions about where to focus their efforts and resources. This proactive approach can lead to more sustainable water infrastructure, ensuring that clean water is accessible to all Tanzanians. Furthermore, the insights gained from this predictive modeling can guide future well construction and maintenance practices, ultimately enhancing the resilience and effectiveness of Tanzania's water supply systems.

## 0.3 Data Understanding

Our data sources are :

**Training Set Values**: 1. Description: Contains independent variables about each water well (e.g., type of pump, installation year, location). 2. Usage: Used to train the predictive model.

**Training Set Labels**:

1. Description: Contains the dependent variable (status_group) for each well, indicating its condition (functional, non-functional, needs repair).
2. Usage: Provides target outcomes for training the model

**Test Set Values**:

1. Description: Contains independent variables for wells needing predictions, similar to the training set values but without labels.
2. Usage: The model predicts the condition of these wells.

**Submission Format**:

1. Description: Template for submitting predictions, including well IDs and predicted status_group.
2. Usage: Ensures predictions are submitted in the correct format for evaluation.

```python
[211]: # Importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import geopandas as gpd
from sklearn.cluster import KMeans
import scipy.stats as stats
import prince
import contextily as ctx
from matplotlib.lines import Line2D
# Load the data
train_values = pd.read_csv("C:/Users/Akipkurui/Desktop/Moringa- Data Science/
 ↪Phase 3/Project/data/training set values.csv")
```

```python
train_labels = pd.read_csv("C:/Users/Akipkurui/Desktop/Moringa- Data Science/
 ↪Phase 3/Project/data/training set labels.csv")
test_values = pd.read_csv("C:/Users/Akipkurui/Desktop/Moringa- Data Science/
 ↪Phase 3/Project/data/test set values.csv")
```

[212]:
```python
# Display the first few rows of the datasets
train_values.head()
```

[212]:

| | id | amount_tsh | date_recorded | funder | gps_height | installer \ |
|---|---|---|---|---|---|---|
| 0 | 69572 | 6000.0 | 14/03/2011 | Roman | 1390 | Roman |
| 1 | 8776 | 0.0 | 06/03/2013 | Grumeti | 1399 | GRUMETI |
| 2 | 34310 | 25.0 | 25/02/2013 | Lottery Club | 686 | World vision |
| 3 | 67743 | 0.0 | 28/01/2013 | Unicef | 263 | UNICEF |
| 4 | 19728 | 0.0 | 13/07/2011 | Action In A | 0 | Artisan |

| | longitude | latitude | wpt_name | num_private | … | payment_type \ |
|---|---|---|---|---|---|---|
| 0 | 34.938093 | -9.856322 | none | 0 | … | annually |
| 1 | 34.698766 | -2.147466 | Zahanati | 0 | … | never pay |
| 2 | 37.460664 | -3.821329 | Kwa Mahundi | 0 | … | per bucket |
| 3 | 38.486161 | -11.155298 | Zahanati Ya Nanyumbu | 0 | … | never pay |
| 4 | 31.130847 | -1.825359 | Shuleni | 0 | … | never pay |

| | water_quality | quality_group | quantity | quantity_group \ |
|---|---|---|---|---|
| 0 | soft | good | enough | enough |
| 1 | soft | good | insufficient | insufficient |
| 2 | soft | good | enough | enough |
| 3 | soft | good | dry | dry |
| 4 | soft | good | seasonal | seasonal |

| | source | source_type | source_class \ |
|---|---|---|---|
| 0 | spring | spring | groundwater |
| 1 | rainwater harvesting | rainwater harvesting | surface |
| 2 | dam | dam | surface |
| 3 | machine dbh | borehole | groundwater |
| 4 | rainwater harvesting | rainwater harvesting | surface |

| | waterpoint_type | waterpoint_type_group |
|---|---|---|
| 0 | communal standpipe | communal standpipe |
| 1 | communal standpipe | communal standpipe |
| 2 | communal standpipe multiple | communal standpipe |
| 3 | communal standpipe multiple | communal standpipe |
| 4 | communal standpipe | communal standpipe |

[5 rows x 40 columns]

[213]:
```python
train_labels.head()
```

```
[213]:       id    status_group
       0  69572        functional
       1   8776        functional
       2  34310        functional
       3  67743    non functional
       4  19728        functional
```

```
[214]: test_values.head()
```

```
[214]:       id  amount_tsh date_recorded                 funder  gps_height  \
       0  50785         0.0    04/02/2013                   Dmdd        1996
       1  51630         0.0    04/02/2013  Government Of Tanzania       1569
       2  17168         0.0    01/02/2013                    NaN        1567
       3  45559         0.0    22/01/2013             Finn Water         267
       4  49871       500.0    27/03/2013                 Bruder        1260

          installer  longitude   latitude                 wpt_name  num_private  \
       0       DMDD  35.290799  -4.059696  Dinamu Secondary School            0
       1        DWE  36.656709  -3.309214                  Kimnyak            0
       2        NaN  34.767863  -5.004344           Puma Secondary            0
       3  FINN WATER  38.058046  -9.418672           Kwa Mzee Pange            0
       4     BRUDER  35.006123 -10.950412          Kwa Mzee Turuka            0

          … payment_type water_quality quality_group       quantity  quantity_group  \
       0  …    never pay          soft          good       seasonal        seasonal
       1  …    never pay          soft          good  insufficient    insufficient
       2  …    never pay          soft          good  insufficient    insufficient
       3  …     unknown          soft          good           dry             dry
       4  …     monthly          soft          good        enough          enough

                       source          source_type  source_class  \
       0  rainwater harvesting  rainwater harvesting       surface
       1                spring                spring   groundwater
       2  rainwater harvesting  rainwater harvesting       surface
       3          shallow well          shallow well   groundwater
       4                spring                spring   groundwater

              waterpoint_type waterpoint_type_group
       0                other                 other
       1  communal standpipe    communal standpipe
       2                other                 other
       3                other                 other
       4  communal standpipe    communal standpipe

       [5 rows x 40 columns]
```

[215]: *#Checking the structure of our datasets*
`train_values.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  int64
 1   amount_tsh           59400 non-null  float64
 2   date_recorded        59400 non-null  object
 3   funder               55763 non-null  object
 4   gps_height           59400 non-null  int64
 5   installer            55745 non-null  object
 6   longitude            59400 non-null  float64
 7   latitude             59400 non-null  float64
 8   wpt_name             59398 non-null  object
 9   num_private          59400 non-null  int64
 10  basin                59400 non-null  object
 11  subvillage           59029 non-null  object
 12  region               59400 non-null  object
 13  region_code          59400 non-null  int64
 14  district_code        59400 non-null  int64
 15  lga                  59400 non-null  object
 16  ward                 59400 non-null  object
 17  population           59400 non-null  int64
 18  public_meeting       56066 non-null  object
 19  recorded_by          59400 non-null  object
 20  scheme_management    55522 non-null  object
 21  scheme_name          30590 non-null  object
 22  permit               56344 non-null  object
 23  construction_year    59400 non-null  int64
 24  extraction_type      59400 non-null  object
 25  extraction_type_group 59400 non-null object
 26  extraction_type_class 59400 non-null object
 27  management           59400 non-null  object
 28  management_group     59400 non-null  object
 29  payment              59400 non-null  object
 30  payment_type         59400 non-null  object
 31  water_quality        59400 non-null  object
 32  quality_group        59400 non-null  object
 33  quantity             59400 non-null  object
 34  quantity_group       59400 non-null  object
 35  source               59400 non-null  object
 36  source_type          59400 non-null  object
 37  source_class         59400 non-null  object
 38  waterpoint_type      59400 non-null  object
```

```
 39   waterpoint_type_group   59400 non-null   object
dtypes: float64(3), int64(7), object(30)
memory usage: 18.1+ MB
```

[216]: `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
```

[217]: `test_values.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14850 entries, 0 to 14849
Data columns (total 40 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   id                 14850 non-null  int64
 1   amount_tsh         14850 non-null  float64
 2   date_recorded      14850 non-null  object
 3   funder             13980 non-null  object
 4   gps_height         14850 non-null  int64
 5   installer          13973 non-null  object
 6   longitude          14850 non-null  float64
 7   latitude           14850 non-null  float64
 8   wpt_name           14850 non-null  object
 9   num_private        14850 non-null  int64
 10  basin              14850 non-null  object
 11  subvillage         14751 non-null  object
 12  region             14850 non-null  object
 13  region_code        14850 non-null  int64
 14  district_code      14850 non-null  int64
 15  lga                14850 non-null  object
 16  ward               14850 non-null  object
 17  population         14850 non-null  int64
 18  public_meeting     14029 non-null  object
 19  recorded_by        14850 non-null  object
 20  scheme_management  13881 non-null  object
 21  scheme_name        7608 non-null   object
 22  permit             14113 non-null  object
 23  construction_year  14850 non-null  int64
 24  extraction_type    14850 non-null  object
```

```
25  extraction_type_group  14850 non-null  object
26  extraction_type_class  14850 non-null  object
27  management             14850 non-null  object
28  management_group       14850 non-null  object
29  payment                14850 non-null  object
30  payment_type           14850 non-null  object
31  water_quality          14850 non-null  object
32  quality_group          14850 non-null  object
33  quantity               14850 non-null  object
34  quantity_group         14850 non-null  object
35  source                 14850 non-null  object
36  source_type            14850 non-null  object
37  source_class           14850 non-null  object
38  waterpoint_type        14850 non-null  object
39  waterpoint_type_group  14850 non-null  object
dtypes: float64(3), int64(7), object(30)
memory usage: 4.5+ MB
```

[218]: `# Getting a summary of statitiscal measures about the data set`
`train_values.describe()`

[218]:
```
                 id      amount_tsh      gps_height       longitude        latitude  \
count  59400.000000    59400.000000    59400.000000    59400.000000    5.940000e+04
mean   37115.131768      317.650385      668.297239       34.077427   -5.706033e+00
std    21453.128371     2997.574558      693.116350        6.567432    2.946019e+00
min        0.000000        0.000000      -90.000000        0.000000   -1.164944e+01
25%    18519.750000        0.000000        0.000000       33.090347   -8.540621e+00
50%    37061.500000        0.000000      369.000000       34.908743   -5.021597e+00
75%    55656.500000       20.000000     1319.250000       37.178387   -3.326156e+00
max    74247.000000   350000.000000     2770.000000       40.345193   -2.000000e-08


        num_private    region_code    district_code      population  \
count  59400.000000   59400.000000     59400.000000    59400.000000
mean       0.474141      15.297003         5.629747      179.909983
std       12.236230      17.587406         9.633649      471.482176
min        0.000000       1.000000         0.000000        0.000000
25%        0.000000       5.000000         2.000000        0.000000
50%        0.000000      12.000000         3.000000       25.000000
75%        0.000000      17.000000         5.000000      215.000000
max     1776.000000      99.000000        80.000000    30500.000000


        construction_year
count        59400.000000
mean          1300.652475
std            951.620547
min              0.000000
25%              0.000000
```

```
50%            1986.000000
75%            2004.000000
max            2013.000000
```

[219]: `train_labels.describe()`

[219]:
```
                    id
count   59400.000000
mean    37115.131768
std     21453.128371
min         0.000000
25%     18519.750000
50%     37061.500000
75%     55656.500000
max     74247.000000
```

[220]: `test_values.describe()`

[220]:
```
                 id      amount_tsh      gps_height      longitude        latitude  \
count  14850.000000   14850.000000   14850.000000   14850.000000   1.485000e+04
mean   37161.972929     322.826983     655.147609      34.061605  -5.684724e+00
std    21359.364833    2510.968644     691.261185       6.593034   2.940803e+00
min       10.000000       0.000000     -57.000000       0.000000  -1.156459e+01
25%    18727.000000       0.000000       0.000000      33.069455  -8.443970e+00
50%    37361.500000       0.000000     344.000000      34.901215  -5.049750e+00
75%    55799.750000      25.000000    1308.000000      37.196594  -3.320594e+00
max    74249.000000  200000.000000    2777.000000      40.325016  -2.000000e-08

          num_private    region_code   district_code     population  \
count    14850.000000   14850.000000    14850.000000   14850.000000
mean         0.415084      15.139057        5.626397     184.114209
std          8.167910      17.191329        9.673842     469.499332
min          0.000000       1.000000        0.000000       0.000000
25%          0.000000       5.000000        2.000000       0.000000
50%          0.000000      12.000000        3.000000      20.000000
75%          0.000000      17.000000        5.000000     220.000000
max        669.000000      99.000000       80.000000   11469.000000

          construction_year
count          14850.000000
mean            1289.708350
std              955.241087
min                0.000000
25%                0.000000
50%             1986.000000
75%             2004.000000
max             2013.000000
```

## 0.4 Data Preparation

### 0.4.1 Merging Datasets

```python
[221]: # Merging train set values and train set labels
       train_data= (pd.merge(train_values,train_labels,on="id",how="inner"))
       train_data
```

```
[221]:            id  amount_tsh date_recorded           funder  gps_height  \
       0       69572      6000.0    14/03/2011           Roman        1390
       1        8776         0.0    06/03/2013         Grumeti        1399
       2       34310        25.0    25/02/2013     Lottery Club         686
       3       67743         0.0    28/01/2013          Unicef         263
       4       19728         0.0    13/07/2011     Action In A           0
       ...       ...         ...           ...             ...         ...
       59395   60739        10.0    03/05/2013  Germany Republi        1210
       59396   27263      4700.0    07/05/2011    Cefa-njombe        1212
       59397   37057         0.0    11/04/2011             NaN           0
       59398   31282         0.0    08/03/2011           Malec           0
       59399   26348         0.0    23/03/2011      World Bank         191


                 installer  longitude    latitude                wpt_name  num_private  \
       0             Roman  34.938093   -9.856322                    none            0
       1           GRUMETI  34.698766   -2.147466                Zahanati            0
       2      World vision  37.460664   -3.821329             Kwa Mahundi            0
       3            UNICEF  38.486161  -11.155298  Zahanati Ya Nanyumbu            0
       4            Artisan  31.130847   -1.825359                 Shuleni            0
       ...             ...        ...         ...                     ...          ...
       59395           CES  37.169807   -3.253847   Area Three Namba 27            0
       59396          Cefa  35.249991   -9.070629    Kwa Yahona Kuvala            0
       59397           NaN  34.017087   -8.750434                 Mashine            0
       59398          Musa  35.861315   -6.378573                  Mshoro            0
       59399         World  38.104048   -6.747464     Kwa Mzee Lugawa            0


              ... water_quality quality_group      quantity  quantity_group  \
       0       ...          soft          good        enough          enough
       1       ...          soft          good  insufficient    insufficient
       2       ...          soft          good        enough          enough
       3       ...          soft          good           dry             dry
       4       ...          soft          good      seasonal        seasonal
       ...     ...           ...           ...           ...             ...
       59395   ...          soft          good        enough          enough
       59396   ...          soft          good        enough          enough
       59397   ...      fluoride      fluoride        enough          enough
       59398   ...          soft          good  insufficient    insufficient
       59399   ...         salty         salty        enough          enough


                          source           source_type source_class  \
```

```
0                    spring                 spring  groundwater
1      rainwater harvesting   rainwater harvesting      surface
2                       dam                    dam      surface
3               machine dbh               borehole  groundwater
4      rainwater harvesting   rainwater harvesting      surface
...                     ...                    ...          ...
59395                 spring                 spring  groundwater
59396                  river             river/lake      surface
59397            machine dbh               borehole  groundwater
59398           shallow well           shallow well  groundwater
59399           shallow well           shallow well  groundwater

                 waterpoint_type waterpoint_type_group    status_group
0              communal standpipe     communal standpipe      functional
1              communal standpipe     communal standpipe      functional
2      communal standpipe multiple     communal standpipe      functional
3      communal standpipe multiple     communal standpipe  non functional
4              communal standpipe     communal standpipe      functional
...                          ...                    ...             ...
59395          communal standpipe     communal standpipe      functional
59396          communal standpipe     communal standpipe      functional
59397                   hand pump              hand pump      functional
59398                   hand pump              hand pump      functional
59399                   hand pump              hand pump      functional

[59400 rows x 41 columns]
```

### 0.4.2  Duplicates removal

```
[222]: # Checking for duplicates
       duplicates = train_data[train_data.duplicated()]
       duplicates
```

```
[222]: Empty DataFrame
       Columns: [id, amount_tsh, date_recorded, funder, gps_height, installer,
       longitude, latitude, wpt_name, num_private, basin, subvillage, region,
       region_code, district_code, lga, ward, population, public_meeting, recorded_by,
       scheme_management, scheme_name, permit, construction_year, extraction_type,
       extraction_type_group, extraction_type_class, management, management_group,
       payment, payment_type, water_quality, quality_group, quantity, quantity_group,
       source, source_type, source_class, waterpoint_type, waterpoint_type_group,
       status_group]
       Index: []

       [0 rows x 41 columns]
```

There are no duplicate rows in this data set

### 0.4.3 Dealing with null values

```python
# Checking for columns with null values
print(train_data.isna().mean()*100)
```

```
id                        0.000000
amount_tsh                0.000000
date_recorded             0.000000
funder                    6.122896
gps_height                0.000000
installer                 6.153199
longitude                 0.000000
latitude                  0.000000
wpt_name                  0.003367
num_private               0.000000
basin                     0.000000
subvillage                0.624579
region                    0.000000
region_code               0.000000
district_code             0.000000
lga                       0.000000
ward                      0.000000
population                0.000000
public_meeting            5.612795
recorded_by               0.000000
scheme_management         6.528620
scheme_name              48.501684
permit                    5.144781
construction_year         0.000000
extraction_type           0.000000
extraction_type_group     0.000000
extraction_type_class     0.000000
management                0.000000
management_group          0.000000
payment                   0.000000
payment_type              0.000000
water_quality             0.000000
quality_group             0.000000
quantity                  0.000000
quantity_group            0.000000
source                    0.000000
source_type               0.000000
source_class              0.000000
waterpoint_type           0.000000
waterpoint_type_group     0.000000
status_group              0.000000
dtype: float64
```

### 0.4.4 Checking for similar columns

**Scheme_Management, Management group and Management**

```
[224]: # Grouping the management columns to check similarity
       train_data['management'].value_counts()
```

```
[224]: management
       vwc                40507
       wug                 6515
       water board         2933
       wua                 2535
       private operator    1971
       parastatal          1768
       water authority      904
       other                844
       company              685
       unknown              561
       other - school        99
       trust                 78
       Name: count, dtype: int64
```

```
[225]: train_data['management_group'].value_counts()
```

```
[225]: management_group
       user-group    52490
       commercial     3638
       parastatal     1768
       other           943
       unknown         561
       Name: count, dtype: int64
```

```
[226]: train_data['scheme_management'].value_counts()
```

```
[226]: scheme_management
       VWC                36793
       WUG                 5206
       Water authority     3153
       WUA                 2883
       Water Board         2748
       Parastatal          1680
       Private operator    1063
       Company             1061
       Other                766
       SWC                   97
       Trust                 72
       Name: count, dtype: int64
```

```
[227]: train_data.groupby(['management_group','management']).count()
```

[227]:

|  |  | id | amount_tsh | date_recorded | funder |
|---|---|---|---|---|---|
| management_group | management |  |  |  |  |
| commercial | company | 685 | 685 | 685 | 663 |
|  | private operator | 1971 | 1971 | 1971 | 1957 |
|  | trust | 78 | 78 | 78 | 78 |
|  | water authority | 904 | 904 | 904 | 836 |
| other | other | 844 | 844 | 844 | 837 |
|  | other - school | 99 | 99 | 99 | 99 |
| parastatal | parastatal | 1768 | 1768 | 1768 | 1624 |
| unknown | unknown | 561 | 561 | 561 | 533 |
| user-group | vwc | 40507 | 40507 | 40507 | 37630 |
|  | water board | 2933 | 2933 | 2933 | 2715 |
|  | wua | 2535 | 2535 | 2535 | 2308 |
|  | wug | 6515 | 6515 | 6515 | 6483 |

|  |  | gps_height | installer | longitude | latitude |
|---|---|---|---|---|---|
| management_group | management |  |  |  |  |
| commercial | company | 685 | 663 | 685 | 685 |
|  | private operator | 1971 | 1959 | 1971 | 1971 |
|  | trust | 78 | 78 | 78 | 78 |
|  | water authority | 904 | 836 | 904 | 904 |
| other | other | 844 | 831 | 844 | 844 |
|  | other - school | 99 | 99 | 99 | 99 |
| parastatal | parastatal | 1768 | 1626 | 1768 | 1768 |
| unknown | unknown | 561 | 527 | 561 | 561 |
| user-group | vwc | 40507 | 37630 | 40507 | 40507 |
|  | water board | 2933 | 2714 | 2933 | 2933 |
|  | wua | 2535 | 2309 | 2535 | 2535 |
|  | wug | 6515 | 6473 | 6515 | 6515 |

|  |  | wpt_name | num_private | … | water_quality |
|---|---|---|---|---|---|
| management_group | management |  |  | … |  |
| commercial | company | 685 | 685 | … | 685 |
|  | private operator | 1971 | 1971 | … | 1971 |
|  | trust | 78 | 78 | … | 78 |
|  | water authority | 904 | 904 | … | 904 |
| other | other | 844 | 844 | … | 844 |
|  | other - school | 99 | 99 | … | 99 |
| parastatal | parastatal | 1768 | 1768 | … | 1768 |
| unknown | unknown | 561 | 561 | … | 561 |
| user-group | vwc | 40507 | 40507 | … | 40507 |
|  | water board | 2932 | 2933 | … | 2933 |
|  | wua | 2535 | 2535 | … | 2535 |
|  | wug | 6514 | 6515 | … | 6515 |

|  |  | quality_group | quantity | quantity_group \ |
| --- | --- | --- | --- | --- |
| management_group | management |  |  |  |
| commercial | company | 685 | 685 | 685 |
|  | private operator | 1971 | 1971 | 1971 |
|  | trust | 78 | 78 | 78 |
|  | water authority | 904 | 904 | 904 |
| other | other | 844 | 844 | 844 |
|  | other - school | 99 | 99 | 99 |
| parastatal | parastatal | 1768 | 1768 | 1768 |
| unknown | unknown | 561 | 561 | 561 |
| user-group | vwc | 40507 | 40507 | 40507 |
|  | water board | 2933 | 2933 | 2933 |
|  | wua | 2535 | 2535 | 2535 |
|  | wug | 6515 | 6515 | 6515 |

|  |  | source | source_type | source_class \ |
| --- | --- | --- | --- | --- |
| management_group | management |  |  |  |
| commercial | company | 685 | 685 | 685 |
|  | private operator | 1971 | 1971 | 1971 |
|  | trust | 78 | 78 | 78 |
|  | water authority | 904 | 904 | 904 |
| other | other | 844 | 844 | 844 |
|  | other - school | 99 | 99 | 99 |
| parastatal | parastatal | 1768 | 1768 | 1768 |
| unknown | unknown | 561 | 561 | 561 |
| user-group | vwc | 40507 | 40507 | 40507 |
|  | water board | 2933 | 2933 | 2933 |
|  | wua | 2535 | 2535 | 2535 |
|  | wug | 6515 | 6515 | 6515 |

|  |  | waterpoint_type | waterpoint_type_group \ |
| --- | --- | --- | --- |
| management_group | management |  |  |
| commercial | company | 685 | 685 |
|  | private operator | 1971 | 1971 |
|  | trust | 78 | 78 |
|  | water authority | 904 | 904 |
| other | other | 844 | 844 |
|  | other - school | 99 | 99 |
| parastatal | parastatal | 1768 | 1768 |
| unknown | unknown | 561 | 561 |
| user-group | vwc | 40507 | 40507 |
|  | water board | 2933 | 2933 |
|  | wua | 2535 | 2535 |
|  | wug | 6515 | 6515 |

|  |  | status_group |
| --- | --- | --- |
| management_group | management |  |

```
commercial       company                    685
                 private operator          1971
                 trust                       78
                 water authority            904
other            other                      844
                 other - school              99
parastatal       parastatal                1768
unknown          unknown                    561
user-group       vwc                      40507
                 water board               2933
                 wua                       2535
                 wug                       6515

[12 rows x 39 columns]
```

From above data scheme_management and management columns have the same values however scheme_management has null values hence we wll drop scheme_management. management column is a subset of management-group hence similarity. Since Management column is more detailed, we will drop management_group and retain management column

**extraction_type, extraction_type_group and extraction_type_class**

```
[228]: train_data['extraction_type'].value_counts()
```

```
[228]: extraction_type
       gravity                    26780
       nira/tanira                 8154
       other                       6430
       submersible                 4764
       swn 80                      3670
       mono                        2865
       india mark ii               2400
       afridev                     1770
       ksb                         1415
       other - rope pump            451
       other - swn 81               229
       windmill                     117
       india mark iii                98
       cemo                          90
       other - play pump             85
       walimi                        48
       climax                        32
       other - mkulima/shinyanga      2
       Name: count, dtype: int64
```

```
[229]: train_data['extraction_type_group'].value_counts()
```

```
[229]: extraction_type_group
       gravity          26780
       nira/tanira       8154
       other             6430
       submersible       6179
       swn 80            3670
       mono              2865
       india mark ii     2400
       afridev           1770
       rope pump          451
       other handpump     364
       other motorpump    122
       wind-powered       117
       india mark iii      98
       Name: count, dtype: int64
```

```
[230]: train_data['extraction_type_class'].value_counts()
```

```
[230]: extraction_type_class
       gravity        26780
       handpump       16456
       other           6430
       submersible     6179
       motorpump       2987
       rope pump        451
       wind-powered     117
       Name: count, dtype: int64
```

```
[231]: train_data.groupby(['extraction_type_class','extraction_type_group']).count()
```

| [231]: | | id | amount_tsh | date_recorded \ |
|---|---|---|---|---|
| extraction_type_class | extraction_type_group | | | |
| gravity | gravity | 26780 | 26780 | 26780 |
| handpump | afridev | 1770 | 1770 | 1770 |
| | india mark ii | 2400 | 2400 | 2400 |
| | india mark iii | 98 | 98 | 98 |
| | nira/tanira | 8154 | 8154 | 8154 |
| | other handpump | 364 | 364 | 364 |
| | swn 80 | 3670 | 3670 | 3670 |
| motorpump | mono | 2865 | 2865 | 2865 |
| | other motorpump | 122 | 122 | 122 |
| other | other | 6430 | 6430 | 6430 |
| rope pump | rope pump | 451 | 451 | 451 |
| submersible | submersible | 6179 | 6179 | 6179 |
| wind-powered | wind-powered | 117 | 117 | 117 |

```
                         funder  gps_height  installer  \
```

| extraction_type_class | extraction_type_group | | | |
|---|---|---|---|---|
| gravity | gravity | 24704 | 26780 | 24714 |
| handpump | afridev | 1668 | 1770 | 1665 |
| | india mark ii | 2358 | 2400 | 2358 |
| | india mark iii | 98 | 98 | 98 |
| | nira/tanira | 7899 | 8154 | 7885 |
| | other handpump | 353 | 364 | 354 |
| | swn 80 | 3595 | 3670 | 3593 |
| motorpump | mono | 2577 | 2865 | 2578 |
| | other motorpump | 122 | 122 | 122 |
| other | other | 6010 | 6430 | 6002 |
| rope pump | rope pump | 448 | 451 | 448 |
| submersible | submersible | 5819 | 6179 | 5816 |
| wind-powered | wind-powered | 112 | 117 | 112 |

| | | longitude | latitude | wpt_name \ |
|---|---|---|---|---|
| extraction_type_class | extraction_type_group | | | |
| gravity | gravity | 26780 | 26780 | 26779 |
| handpump | afridev | 1770 | 1770 | 1770 |
| | india mark ii | 2400 | 2400 | 2400 |
| | india mark iii | 98 | 98 | 98 |
| | nira/tanira | 8154 | 8154 | 8154 |
| | other handpump | 364 | 364 | 364 |
| | swn 80 | 3670 | 3670 | 3670 |
| motorpump | mono | 2865 | 2865 | 2865 |
| | other motorpump | 122 | 122 | 122 |
| other | other | 6430 | 6430 | 6429 |
| rope pump | rope pump | 451 | 451 | 451 |
| submersible | submersible | 6179 | 6179 | 6179 |
| wind-powered | wind-powered | 117 | 117 | 117 |

| | | num_private | … | water_quality \ |
|---|---|---|---|---|
| extraction_type_class | extraction_type_group | | … | |
| gravity | gravity | 26780 | … | 26780 |
| handpump | afridev | 1770 | … | 1770 |
| | india mark ii | 2400 | … | 2400 |
| | india mark iii | 98 | … | 98 |
| | nira/tanira | 8154 | … | 8154 |
| | other handpump | 364 | … | 364 |
| | swn 80 | 3670 | … | 3670 |
| motorpump | mono | 2865 | … | 2865 |
| | other motorpump | 122 | … | 122 |
| other | other | 6430 | … | 6430 |
| rope pump | rope pump | 451 | … | 451 |
| submersible | submersible | 6179 | … | 6179 |
| wind-powered | wind-powered | 117 | … | 117 |

| extraction_type_class | extraction_type_group | quality_group | quantity \ |
|---|---|---|---|
| gravity | gravity | 26780 | 26780 |
| handpump | afridev | 1770 | 1770 |
| | india mark ii | 2400 | 2400 |
| | india mark iii | 98 | 98 |
| | nira/tanira | 8154 | 8154 |
| | other handpump | 364 | 364 |
| | swn 80 | 3670 | 3670 |
| motorpump | mono | 2865 | 2865 |
| | other motorpump | 122 | 122 |
| other | other | 6430 | 6430 |
| rope pump | rope pump | 451 | 451 |
| submersible | submersible | 6179 | 6179 |
| wind-powered | wind-powered | 117 | 117 |

| extraction_type_class | extraction_type_group | quantity_group | source \ |
|---|---|---|---|
| gravity | gravity | 26780 | 26780 |
| handpump | afridev | 1770 | 1770 |
| | india mark ii | 2400 | 2400 |
| | india mark iii | 98 | 98 |
| | nira/tanira | 8154 | 8154 |
| | other handpump | 364 | 364 |
| | swn 80 | 3670 | 3670 |
| motorpump | mono | 2865 | 2865 |
| | other motorpump | 122 | 122 |
| other | other | 6430 | 6430 |
| rope pump | rope pump | 451 | 451 |
| submersible | submersible | 6179 | 6179 |
| wind-powered | wind-powered | 117 | 117 |

| extraction_type_class | extraction_type_group | source_type | source_class \ |
|---|---|---|---|
| gravity | gravity | 26780 | 26780 |
| handpump | afridev | 1770 | 1770 |
| | india mark ii | 2400 | 2400 |
| | india mark iii | 98 | 98 |
| | nira/tanira | 8154 | 8154 |
| | other handpump | 364 | 364 |
| | swn 80 | 3670 | 3670 |
| motorpump | mono | 2865 | 2865 |
| | other motorpump | 122 | 122 |
| other | other | 6430 | 6430 |
| rope pump | rope pump | 451 | 451 |
| submersible | submersible | 6179 | 6179 |
| wind-powered | wind-powered | 117 | 117 |

|  |  | waterpoint_type \ |
| --- | --- | --- |
| extraction_type_class | extraction_type_group |  |
| gravity | gravity | 26780 |
| handpump | afridev | 1770 |
|  | india mark ii | 2400 |
|  | india mark iii | 98 |
|  | nira/tanira | 8154 |
|  | other handpump | 364 |
|  | swn 80 | 3670 |
| motorpump | mono | 2865 |
|  | other motorpump | 122 |
| other | other | 6430 |
| rope pump | rope pump | 451 |
| submersible | submersible | 6179 |
| wind-powered | wind-powered | 117 |

|  |  | waterpoint_type_group \ |
| --- | --- | --- |
| extraction_type_class | extraction_type_group |  |
| gravity | gravity | 26780 |
| handpump | afridev | 1770 |
|  | india mark ii | 2400 |
|  | india mark iii | 98 |
|  | nira/tanira | 8154 |
|  | other handpump | 364 |
|  | swn 80 | 3670 |
| motorpump | mono | 2865 |
|  | other motorpump | 122 |
| other | other | 6430 |
| rope pump | rope pump | 451 |
| submersible | submersible | 6179 |
| wind-powered | wind-powered | 117 |

|  |  | status_group |
| --- | --- | --- |
| extraction_type_class | extraction_type_group |  |
| gravity | gravity | 26780 |
| handpump | afridev | 1770 |
|  | india mark ii | 2400 |
|  | india mark iii | 98 |
|  | nira/tanira | 8154 |
|  | other handpump | 364 |
|  | swn 80 | 3670 |
| motorpump | mono | 2865 |
|  | other motorpump | 122 |
| other | other | 6430 |
| rope pump | rope pump | 451 |
| submersible | submersible | 6179 |

```
wind-powered           wind-powered                    117
```

```
[13 rows x 39 columns]
```

It is evident that this 3 columns contain the same data. extraction_type and extraction_type_group contain the same data however extraction_type_group appears to be more compact. extraction_type_group appears to be a subset of extrcation_type_class. We will drop extraction_type and extraction_type_class and retain extraction_type_group since it is more compact and has more details.

**payment and payment_type**

```
[232]: train_data['payment'].value_counts()
```

```
[232]: payment
       never pay                25348
       pay per bucket            8985
       pay monthly               8300
       unknown                   8157
       pay when scheme fails     3914
       pay annually              3642
       other                     1054
       Name: count, dtype: int64
```

```
[233]: train_data['payment_type'].value_counts()
```

```
[233]: payment_type
       never pay      25348
       per bucket      8985
       monthly         8300
       unknown         8157
       on failure      3914
       annually        3642
       other           1054
       Name: count, dtype: int64
```

**This 2 columns are similar we decided to drop 1 i.e payment**

**quantity and quantity-group**

```
[234]: train_data['quantity'].value_counts()
```

```
[234]: quantity
       enough          33186
       insufficient    15129
       dry              6246
       seasonal         4050
       unknown           789
       Name: count, dtype: int64
```

```
[235]: train_data['quantity_group'].value_counts()
```

```
[235]: quantity_group
       enough          33186
       insufficient    15129
       dry              6246
       seasonal         4050
       unknown           789
       Name: count, dtype: int64
```

**This 2 columns are similar we decided to drop 1 i.e quantity_group**

**water_quality and quality_group**

```
[236]: train_data['water_quality'].value_counts()
```

```
[236]: water_quality
       soft                50818
       salty                4856
       unknown              1876
       milky                 804
       coloured              490
       salty abandoned       339
       fluoride              200
       fluoride abandoned     17
       Name: count, dtype: int64
```

```
[237]: train_data['quality_group'].value_counts()
```

```
[237]: quality_group
       good       50818
       salty       5195
       unknown     1876
       milky        804
       colored      490
       fluoride     217
       Name: count, dtype: int64
```

**The 2 columns are similar however water_quality column has more details hence we will drop quality_group**

**source, source_type and source_class**

```
[238]: train_data['source'].value_counts()
```

```
[238]: source
       spring          17021
       shallow well    16824
       machine dbh     11075
```

```
river                    9612
rainwater harvesting     2295
hand dtw                  874
lake                      765
dam                       656
other                     212
unknown                    66
Name: count, dtype: int64
```

[239]: `train_data['source_type'].value_counts()`

[239]:
```
source_type
spring                  17021
shallow well            16824
borehole                11949
river/lake              10377
rainwater harvesting     2295
dam                       656
other                     278
Name: count, dtype: int64
```

[240]: `train_data['source_class'].value_counts()`

[240]:
```
source_class
groundwater    45794
surface        13328
unknown          278
Name: count, dtype: int64
```

[241]: `train_data.groupby(['source_class','source_type']).count()`

[241]:

| source_class | source_type | id | amount_tsh | date_recorded | funder |
|---|---|---|---|---|---|
| groundwater | borehole | 11949 | 11949 | 11949 | 11119 |
|  | shallow well | 16824 | 16824 | 16824 | 16301 |
|  | spring | 17021 | 17021 | 17021 | 15870 |
| surface | dam | 656 | 656 | 656 | 647 |
|  | rainwater harvesting | 2295 | 2295 | 2295 | 2099 |
|  | river/lake | 10377 | 10377 | 10377 | 9478 |
| unknown | other | 278 | 278 | 278 | 249 |

| source_class | source_type | gps_height | installer | longitude | latitude |
|---|---|---|---|---|---|
| groundwater | borehole | 11949 | 11114 | 11949 | 11949 |
|  | shallow well | 16824 | 16286 | 16824 | 16824 |
|  | spring | 17021 | 15870 | 17021 | 17021 |
| surface | dam | 656 | 646 | 656 | 656 |

|  |  | | | | |
|---|---|---|---|---|---|
|  | rainwater harvesting | 2295 | 2096 | 2295 | 2295 |
|  | river/lake | 10377 | 9483 | 10377 | 10377 |
| unknown | other | 278 | 250 | 278 | 278 |

| source_class | source_type | wpt_name | num_private | … | payment |
|---|---|---|---|---|---|
|  |  |  |  | … |  |
| groundwater | borehole | 11948 | 11949 | … | 11949 |
|  | shallow well | 16824 | 16824 | … | 16824 |
|  | spring | 17020 | 17021 | … | 17021 |
| surface | dam | 656 | 656 | … | 656 |
|  | rainwater harvesting | 2295 | 2295 | … | 2295 |
|  | river/lake | 10377 | 10377 | … | 10377 |
| unknown | other | 278 | 278 | … | 278 |

| source_class | source_type | payment_type | water_quality | quality_group |
|---|---|---|---|---|
| groundwater | borehole | 11949 | 11949 | 11949 |
|  | shallow well | 16824 | 16824 | 16824 |
|  | spring | 17021 | 17021 | 17021 |
| surface | dam | 656 | 656 | 656 |
|  | rainwater harvesting | 2295 | 2295 | 2295 |
|  | river/lake | 10377 | 10377 | 10377 |
| unknown | other | 278 | 278 | 278 |

| source_class | source_type | quantity | quantity_group | source |
|---|---|---|---|---|
| groundwater | borehole | 11949 | 11949 | 11949 |
|  | shallow well | 16824 | 16824 | 16824 |
|  | spring | 17021 | 17021 | 17021 |
| surface | dam | 656 | 656 | 656 |
|  | rainwater harvesting | 2295 | 2295 | 2295 |
|  | river/lake | 10377 | 10377 | 10377 |
| unknown | other | 278 | 278 | 278 |

| source_class | source_type | waterpoint_type | waterpoint_type_group |
|---|---|---|---|
| groundwater | borehole | 11949 | 11949 |
|  | shallow well | 16824 | 16824 |
|  | spring | 17021 | 17021 |
| surface | dam | 656 | 656 |
|  | rainwater harvesting | 2295 | 2295 |
|  | river/lake | 10377 | 10377 |
| unknown | other | 278 | 278 |

| source_class | source_type | status_group |
|---|---|---|
| groundwater | borehole | 11949 |

```
                shallow well          16824
                spring                17021
    surface     dam                     656
                rainwater harvesting   2295
                river/lake            10377
    unknown     other                   278

[7 rows x 39 columns]
```

The 3 columns are similar. source is more detailed that source_type while source_type is a subset of source_class hence we will remain with source and drop source_type and source_class

**waterpoint_type and waterpoint_type_group**

[242]: `train_data['waterpoint_type'].value_counts()`

```
[242]: waterpoint_type
    communal standpipe           28522
    hand pump                    17488
    other                         6380
    communal standpipe multiple   6103
    improved spring                784
    cattle trough                  116
    dam                              7
    Name: count, dtype: int64
```

[243]: `train_data['waterpoint_type_group'].value_counts()`

```
[243]: waterpoint_type_group
    communal standpipe    34625
    hand pump             17488
    other                  6380
    improved spring         784
    cattle trough           116
    dam                       7
    Name: count, dtype: int64
```

The 2 columns are similar howoever waterpoint_type has more details hence we will drop waterpoint_type_group

### 0.4.5 Drop identical columns

[244]: 
```
train_data1=train_data.
 ↪drop(columns=['management_group','scheme_management','extraction_type_class','extraction_ty
              'payment_type', 'waterpoint_type_group'])
train_data1
```

```
[244]:              id  amount_tsh date_recorded             funder  gps_height  \
       0         69572      6000.0    14/03/2011            Roman        1390
       1          8776         0.0    06/03/2013           Grumeti        1399
       2         34310        25.0    25/02/2013      Lottery Club         686
       3         67743         0.0    28/01/2013            Unicef         263
       4         19728         0.0    13/07/2011         Action In A           0
       ...         ...         ...           ...               ...         ...
       59395     60739        10.0    03/05/2013  Germany Republi        1210
       59396     27263      4700.0    07/05/2011       Cefa-njombe        1212
       59397     37057         0.0    11/04/2011               NaN           0
       59398     31282         0.0    08/03/2011             Malec           0
       59399     26348         0.0    23/03/2011        World Bank         191

                 installer   longitude    latitude               wpt_name  num_private  \
       0            Roman  34.938093   -9.856322                   none            0
       1          GRUMETI  34.698766   -2.147466               Zahanati            0
       2     World vision  37.460664   -3.821329            Kwa Mahundi            0
       3           UNICEF  38.486161  -11.155298  Zahanati Ya Nanyumbu            0
       4          Artisan  31.130847   -1.825359                Shuleni            0
       ...            ...         ...         ...                    ...          ...
       59395          CES  37.169807   -3.253847   Area Three Namba 27            0
       59396         Cefa  35.249991   -9.070629     Kwa Yahona Kuvala            0
       59397          NaN  34.017087   -8.750434                Mashine            0
       59398         Musa  35.861315   -6.378573                 Mshoro            0
       59399        World  38.104048   -6.747464      Kwa Mzee Lugawa            0

             … permit construction_year extraction_type_group    management  \
       0     …  False              1999                gravity           vwc
       1     …   True              2010                gravity           wug
       2     …   True              2009                gravity           vwc
       3     …   True              1986            submersible           vwc
       4     …   True                 0                gravity         other
       ...   …    ...               ...                    ...           ...
       59395 …   True              1999                gravity   water board
       59396 …   True              1996                gravity           vwc
       59397 …  False                 0                 swn 80           vwc
       59398 …   True                 0            nira/tanira           vwc
       59399 …   True              2002            nira/tanira           vwc

                     payment water_quality      quantity  \
       0         pay annually          soft        enough
       1            never pay          soft  insufficient
       2       pay per bucket          soft        enough
       3            never pay          soft           dry
       4            never pay          soft      seasonal
       ...              ...           ...           ...
       59395     pay per bucket        soft        enough
```

```
59396              pay annually          soft          enough
59397               pay monthly      fluoride          enough
59398                  never pay          soft   insufficient
59399   pay when scheme fails         salty          enough


                         source               waterpoint_type   status_group
0                        spring             communal standpipe     functional
1           rainwater harvesting            communal standpipe     functional
2                           dam   communal standpipe multiple     functional
3                   machine dbh   communal standpipe multiple  non functional
4           rainwater harvesting            communal standpipe     functional
...                         ...                           ...            ...
59395                    spring            communal standpipe     functional
59396                     river            communal standpipe     functional
59397               machine dbh                     hand pump     functional
59398               shallow well                    hand pump     functional
59399               shallow well                    hand pump     functional

[59400 rows x 31 columns]
```

**Exporing contruction__year column**

```
[245]: train_data1['construction_year'].value_counts()
```

```
[245]: construction_year
       0        20709
       2010      2645
       2008      2613
       2009      2533
       2000      2091
       2007      1587
       2006      1471
       2003      1286
       2011      1256
       2004      1123
       2012      1084
       2002      1075
       1978      1037
       1995      1014
       2005      1011
       1999       979
       1998       966
       1990       954
       1985       945
       1980       811
       1996       811
       1984       779
```

```
1982    744
1994    738
1972    708
1974    676
1997    644
1992    640
1993    608
2001    540
1988    521
1983    488
1975    437
1986    434
1976    414
1970    411
1991    324
1989    316
1987    302
1981    238
1977    202
1979    192
1973    184
2013    176
1971    145
1960    102
1967     88
1963     85
1968     77
1969     59
1964     40
1962     30
1961     21
1965     19
1966     17
Name: count, dtype: int64
```

[246]: 
```python
train_data1['construction_year'].value_counts().sum()
```

[246]: 59400

**We have replace contruction_year that is balnk with unknown**

[247]: 
```python
# Replace 'Unknown' with NaN
train_data1['construction_year'] = train_data1['construction_year'].
 ↪replace('Unknown', np.nan)

# Convert construction_year to numeric, handling NaN values
```

```python
train_data1['construction_year'] = pd.
  ↪to_numeric(train_data1['construction_year'], errors='coerce')

# Define the function to get the decade in the desired format
def get_decade(year):
    if pd.isna(year):
        return 'Unknown'
    else:
        decade_start = (year // 10) * 10
        return f"{str(int(decade_start))[-2:]}s"

# Apply the function to create the 'decade' column
train_data1['decade'] = train_data1['construction_year'].apply(get_decade)

# Verify the changes
train_data1['decade'].value_counts()
```

[247]: decade
       0s      20709
       00s     15330
       90s      7678
       80s      5578
       10s      5161
       70s      4406
       60s       538
       Name: count, dtype: int64

We have now grouped construction year into decades for easy visualization and interpretation

**recorded_by**

[248]: ```python
train_data1['recorded_by'].value_counts()
```

[248]: recorded_by
       GeoData Consultants Ltd    59400
       Name: count, dtype: int64

[249]: ```python
train_data1.drop(columns=['recorded_by'],inplace=True  )
```

Has only one record hence we will drop this column

### 0.4.6 Correcting errors and spelling mistakes in installer and funder columns

**installer**

[250]: ```python
train_data1['installer'].value_counts()
```

```
[250]: installer
       DWE                  17402
       Government            1825
       RWE                   1206
       Commu                 1060
       DANIDA                1050
                             …
       Wizara  ya maji          1
       TWESS                    1
       Nasan workers            1
       R                        1
       SELEPTA                  1
       Name: count, Length: 2145, dtype: int64
```

```python
[251]: # filling 0 values with unknown
       train_data1['installer'].replace(to_replace = '0', value ='Unknown' ,␣
        ↪inplace=True)
```

```python
[252]: # filling null values with unknown
       train_data1['installer'].fillna(value='Unknown',inplace=True)
```

```python
[253]: # From the most common 100 value counts we realized some spelling mistakes or␣
        ↪different syntax between same categories

       # Replacing the spelling mistakes and collect same categories in same name

       train_data1['installer'].replace(to_replace = ('District Water Department',␣
        ↪'District water depar','Distric Water Department'),
                             value ='District water department' , inplace=True)

       train_data1['installer'].replace(to_replace = ('FinW','Fini water','FINI␣
        ↪WATER'), value ='Fini Water' , inplace=True)
       train_data1['installer'].replace(to_replace = 'JAICA', value ='Jaica' ,␣
        ↪inplace=True)

       train_data1['installer'].replace(to_replace = ('COUN', 'District COUNCIL',␣
        ↪'DISTRICT COUNCIL','District Counci',
                                         'District␣
        ↪Council','Council','Counc','District  Council','Distri'),
                             value ='District council' , inplace=True)

       train_data1['installer'].replace(to_replace = ('RC CHURCH', 'RC Churc',␣
        ↪'RC','RC Ch','RC C', 'RC CH','RC church',
                                         'RC CATHORIC',) , value ='RC Church' ,␣
        ↪inplace=True)
```

```python
train_data1['installer'].replace(to_replace = ('Central Government','Tanzania␣
 ↪Government',
                                    'central government','Cental␣
 ↪Government', 'Cebtral Government',
                                    'Tanzanian Government','Tanzania␣
 ↪government', 'Centra Government' ,
                                    'CENTRAL GOVERNMENT', 'TANZANIAN␣
 ↪GOVERNMENT','Central govt', 'Centr',
                                    'Centra govt') , value ='Central␣
 ↪government' , inplace=True)

train_data1['installer'].replace(to_replace = ('World vision', 'World␣
 ↪Division','World Vision'),
                                        value ='world vision' , inplace=True)

train_data1['installer'].replace(to_replace = ('Unisef','UNICEF'),value␣
 ↪='Unicef' , inplace=True)
train_data1['installer'].replace(to_replace = 'DANID', value ='DANIDA' ,␣
 ↪inplace=True)

train_data1['installer'].replace(to_replace = ('villigers', 'villager',␣
 ↪'Villagers', 'Villa', 'Village', 'Villi',
                                    'Village Council','Village Counil',␣
 ↪'Villages', 'Vill', 'Village community',
                                    'Villaers', 'Village Community',␣
 ↪'Villag','Villege Council', 'Village council',
                                    'Village  Council','Villagerd',␣
 ↪'Villager', 'Village Technician',
                                    'Village Office','Village community␣
 ↪members'),
                                    value ='villagers' , inplace=True)
train_data1['installer'].replace(to_replace␣
 ↪=('Commu','Communit','commu','COMMU', 'COMMUNITY') ,
                                    value ='Community' , inplace=True)

train_data1['installer'].replace(to_replace = ('GOVERNMENT', 'GOVER',␣
 ↪'GOVERNME', 'GOVERM','GOVERN','Gover','Gove',
                                    'Governme','Governmen' ) ,value␣
 ↪='Government' , inplace=True)

train_data1['installer'].replace(to_replace = 'Hesawa' ,value ='HESAWA' ,␣
 ↪inplace=True)
```

```python
[254]: # continue to replacing spellin mistakes and getting together values
train_data1['installer'].replace(to_replace = ('Colonial Government') , value␣
 ↪='Colonial government' , inplace=True)
```

```python
train_data1['installer'].replace(to_replace = ('Government of Misri') , value
 ='Misri Government' , inplace=True)
train_data1['installer'].replace(to_replace = ('Italy government') , value
 ='Italian government' , inplace=True)
train_data1['installer'].replace(to_replace = ('British colonial government') ,
 value ='British government' , inplace=True)
train_data1['installer'].replace(to_replace = ('Concern /government') , value
 ='Concern/Government' , inplace=True)
train_data1['installer'].replace(to_replace = ('Village Government') , value
 ='Village government' , inplace=True)
train_data1['installer'].replace(to_replace = ('Government and Community') ,
 value ='Government /Community' , inplace=True)
train_data1['installer'].replace(to_replace = ('Cetral government /RC') , value
 ='RC church/Central Gover' , inplace=True)
train_data1['installer'].replace(to_replace = ('Government /TCRS','Government/
 TCRS') , value ='TCRS /Government' , inplace=True)
train_data1['installer'].replace(to_replace = ('ADRA /Government') , value
 ='ADRA/Government' , inplace=True)
```

```python
[255]: train_data1['installer'].value_counts().head(20)
```

```
[255]: installer
       DWE                  17402
       Unknown               4435
       Government            2660
       Community             1674
       DANIDA                1602
       HESAWA                1379
       RWE                   1206
       District council      1179
       Central government    1114
       KKKT                   898
       TCRS                   707
       world vision           681
       CES                    610
       Fini Water             593
       RC Church              461
       villagers              408
       LGA                    408
       WEDECO                 397
       TASAF                  396
       Unicef                 332
       Name: count, dtype: int64
```

```python
[256]: # Create a new column 'installer_classified' with default value 'Others'
       train_data1['installer_classified'] = 'Others'
```

```python
# Get the counts of each installer
installer_counts = train_data1['installer'].value_counts()

# Update the 'installer_classified' column for installers with less than 500
 ↪records
for installer, count in installer_counts.items():
    if count > 500:
        train_data1.loc[train_data1['installer'] == installer,
 ↪'installer_classified'] = installer

# Print the first few rows to check the results
train_data1['installer_classified'].value_counts().head(20)
```

[256]: installer_classified
Others               23260
DWE                  17402
Unknown               4435
Government            2660
Community             1674
DANIDA                1602
HESAWA                1379
RWE                   1206
District council      1179
Central government    1114
KKKT                   898
TCRS                   707
world vision           681
CES                    610
Fini Water             593
Name: count, dtype: int64

**funder**

[257]: 
```python
# filling 0 and null values with unknown
train_data1['funder'].fillna(value='Unknown',inplace=True)
train_data1['funder'].replace(to_replace = '0', value ='Unknown' , inplace=True)
train_data1['funder'].value_counts().head(20)
```

[257]: funder
Government Of Tanzania    9084
Unknown                   4418
Danida                    3114
Hesawa                    2202
Rwssp                     1374
World Bank                1349
Kkkt                      1287

```
World Vision              1246
Unicef                    1057
Tasaf                      877
District Council           843
Dhv                        829
Private Individual         826
Dwsp                       811
Norad                      765
Germany Republi            610
Tcrs                       602
Ministry Of Water          590
Water                      583
Dwe                        484
Name: count, dtype: int64
```

[258]:
```python
# Create a new column 'funder_classified' with default value 'Others'
train_data1['funder_classified'] = 'Others'

# Get the counts of each installer
funder_counts = train_data1['funder'].value_counts()

# Update the 'installer_classified' column for installers with less than 500
  ↪records
for funder, count in funder_counts.items():
    if count > 483:
        train_data1.loc[train_data1['funder'] == funder, 'funder_classified'] =
  ↪funder

# Print the first few rows to check the results
train_data1['funder_classified'].value_counts().head(20)
```

[258]:
```
funder_classified
Others                    26449
Government Of Tanzania     9084
Unknown                    4418
Danida                     3114
Hesawa                     2202
Rwssp                      1374
World Bank                 1349
Kkkt                       1287
World Vision               1246
Unicef                     1057
Tasaf                       877
District Council            843
Dhv                         829
Private Individual          826
Dwsp                        811
```

```
Norad                          765
Germany Republi                610
Tcrs                           602
Ministry Of Water              590
Water                          583
Name: count, dtype: int64
```

### 0.4.7 Other columns

We will drop the following columns since they do not have any relationship with functionality of the wells:

1. id
2. wpt_name
3. date_recorded
4. scheme_name
5. region_code

```
[259]: train_data1.
        ↪drop(columns=['wpt_name','scheme_name','id','region_code',"date_recorded"],inplace=True⌴
        ↪)
```

```
[260]: train_data1
```

```
[260]:        amount_tsh            funder  gps_height       installer  longitude  \
       0          6000.0             Roman        1390           Roman  34.938093
       1             0.0           Grumeti        1399         GRUMETI  34.698766
       2            25.0      Lottery Club         686    world vision  37.460664
       3             0.0            Unicef         263          Unicef  38.486161
       4             0.0        Action In A           0         Artisan  31.130847
       ...             ...               ...         ...             ...        ...
       59395        10.0   Germany Republi        1210             CES  37.169807
       59396      4700.0       Cefa-njombe        1212            Cefa  35.249991
       59397         0.0           Unknown           0         Unknown  34.017087
       59398         0.0             Malec           0            Musa  35.861315
       59399         0.0        World Bank         191           World  38.104048

              latitude  num_private                    basin    subvillage  \
       0     -9.856322            0              Lake Nyasa      Mnyusi B
       1     -2.147466            0           Lake Victoria       Nyamara
       2     -3.821329            0                 Pangani       Majengo
       3    -11.155298            0   Ruvuma / Southern Coast  Mahakamani
       4     -1.825359            0           Lake Victoria    Kyanyamisa
       ...          ...          ...                     ...           ...
       59395 -3.253847            0                 Pangani      Kiduruni
       59396 -9.070629            0                  Rufiji      Igumbilo
       59397 -8.750434            0                  Rufiji     Madungulu
       59398 -6.378573            0                  Rufiji        Mwinyi
```

34

```
59399  -6.747464           0            Wami / Ruvu   Kikatanyemba
```

```
           region  …    management              payment water_quality  \
0          Iringa  …           vwc          pay annually          soft
1            Mara  …           wug             never pay          soft
2         Manyara  …           vwc        pay per bucket          soft
3          Mtwara  …           vwc             never pay          soft
4          Kagera  …         other             never pay          soft
…             … …             …                    …             …
59395  Kilimanjaro  …   water board        pay per bucket          soft
59396      Iringa  …           vwc          pay annually          soft
59397       Mbeya  …           vwc           pay monthly      fluoride
59398      Dodoma  …           vwc             never pay          soft
59399    Morogoro  …           vwc  pay when scheme fails         salty
```

```
            quantity              source           waterpoint_type  \
0             enough              spring         communal standpipe
1       insufficient  rainwater harvesting         communal standpipe
2             enough                 dam  communal standpipe multiple
3                dry         machine dbh  communal standpipe multiple
4           seasonal  rainwater harvesting         communal standpipe
…              …                  …                      …
59395         enough              spring         communal standpipe
59396         enough               river         communal standpipe
59397         enough         machine dbh                   hand pump
59398   insufficient        shallow well                   hand pump
59399         enough        shallow well                   hand pump
```

```
           status_group decade installer_classified funder_classified
0              functional     90s               Others            Others
1              functional     10s               Others            Others
2              functional     00s         world vision            Others
3          non functional     80s               Others            Unicef
4              functional      0s               Others            Others
…              …          …                    …                 …
59395          functional     90s                  CES   Germany Republi
59396          functional     90s               Others            Others
59397          functional      0s              Unknown           Unknown
59398          functional      0s               Others            Others
59399          functional     00s               Others        World Bank
```

```
[59400 rows x 28 columns]
```

**amount_tsh**

```
[261]:  train_data1['amount_tsh'].value_counts()
```

```
[261]:  amount_tsh
        0.0        41639
        500.0       3102
        50.0        2472
        1000.0      1488
        20.0        1463
                     …
        6300.0         1
        120000.0       1
        138000.0       1
        350000.0       1
        59.0           1
        Name: count, Length: 98, dtype: int64
```

```
[262]:  train_data1.drop(columns=['amount_tsh'],inplace=True )
```

Most of the records are zero hence we drop the column

```
[263]:  train_data1.drop(columns=['num_private'],inplace=True )
```

```
[264]:  train_data1
```

```
[264]:                    funder  gps_height      installer   longitude   latitude  \
        0                  Roman        1390          Roman   34.938093  -9.856322
        1                Grumeti        1399        GRUMETI   34.698766  -2.147466
        2            Lottery Club         686   world vision   37.460664  -3.821329
        3                 Unicef         263         Unicef   38.486161 -11.155298
        4              Action In A           0        Artisan   31.130847  -1.825359
        …                    …            …              …          …          …
        59395  Germany Republi        1210            CES   37.169807  -3.253847
        59396     Cefa-njombe        1212           Cefa   35.249991  -9.070629
        59397         Unknown           0        Unknown   34.017087  -8.750434
        59398           Malec           0           Musa   35.861315  -6.378573
        59399      World Bank         191          World   38.104048  -6.747464

                                basin    subvillage       region  district_code  \
        0                   Lake Nyasa      Mnyusi B       Iringa              5
        1                Lake Victoria       Nyamara         Mara              2
        2                      Pangani       Majengo      Manyara              4
        3        Ruvuma / Southern Coast   Mahakamani       Mtwara             63
        4                Lake Victoria     Kyanyamisa       Kagera              1
        …                        …             …             …              …
        59395                  Pangani      Kiduruni   Kilimanjaro              5
        59396                   Rufiji      Igumbilo        Iringa              4
        59397                   Rufiji     Madungulu         Mbeya              7
        59398                   Rufiji        Mwinyi        Dodoma              4
        59399             Wami / Ruvu   Kikatanyemba      Morogoro              2
```

```
                   lga  …    management                payment water_quality  \
0               Ludewa  …           vwc            pay annually          soft
1            Serengeti  …           wug              never pay          soft
2            Simanjiro  …           vwc          pay per bucket          soft
3             Nanyumbu  …           vwc              never pay          soft
4              Karagwe  …         other              never pay          soft
…                   …  …           …                     …            …
59395              Hai  …   water board          pay per bucket          soft
59396           Njombe  …           vwc            pay annually          soft
59397          Mbarali  …           vwc             pay monthly      fluoride
59398         Chamwino  …           vwc              never pay          soft
59399   Morogoro Rural  …           vwc  pay when scheme fails         salty

            quantity                source            waterpoint_type  \
0             enough                spring           communal standpipe
1       insufficient   rainwater harvesting           communal standpipe
2             enough                   dam  communal standpipe multiple
3                dry           machine dbh  communal standpipe multiple
4           seasonal   rainwater harvesting           communal standpipe
…                 …                     …                          …
59395         enough                spring           communal standpipe
59396         enough                 river           communal standpipe
59397         enough           machine dbh                    hand pump
59398   insufficient          shallow well                    hand pump
59399         enough          shallow well                    hand pump

            status_group decade installer_classified funder_classified
0             functional    90s               Others            Others
1             functional    10s               Others            Others
2             functional    00s         world vision            Others
3         non functional    80s               Others            Unicef
4             functional     0s               Others            Others
…                     …    …  …                    …                 …
59395         functional    90s                  CES   Germany Republi
59396         functional    90s               Others            Others
59397         functional     0s              Unknown           Unknown
59398         functional     0s               Others            Others
59399         functional    00s               Others        World Bank

[59400 rows x 26 columns]
```

### 0.4.8 Checking remaining null values

```
[265]: print(train_data1.isna().mean()*100)
```

```
funder                  0.000000
```

```
gps_height              0.000000
installer               0.000000
longitude               0.000000
latitude                0.000000
basin                   0.000000
subvillage              0.624579
region                  0.000000
district_code           0.000000
lga                     0.000000
ward                    0.000000
population              0.000000
public_meeting          5.612795
permit                  5.144781
construction_year       0.000000
extraction_type_group   0.000000
management              0.000000
payment                 0.000000
water_quality           0.000000
quantity                0.000000
source                  0.000000
waterpoint_type         0.000000
status_group            0.000000
decade                  0.000000
installer_classified    0.000000
funder_classified       0.000000
dtype: float64
```

### 0.4.9 Remaining columns with null values: subvillage, public_meeting, permit

### 0.4.10 permit

```
[266]: train_data1['permit'].value_counts()
```

```
[266]: permit
       True     38852
       False    17492
       Name: count, dtype: int64
```

```
[267]: train_data1['permit'].fillna(value=True, inplace=True)
```

We have replace null values with True since it is the highest

**public_meeting**

```
[268]: train_data1['public_meeting'].value_counts()
```

```
[268]: public_meeting
       True     51011
       False     5055
```

```
Name: count, dtype: int64
```

[269]: `train_data1['public_meeting'].fillna(value=True, inplace=True)`

**We have replace null values with True since it is the highest**

**subvillage   Drop subvillage since we already have region**

[270]: `train_data1.drop(columns=['subvillage'],inplace=True )`

[271]: `print(train_data1.isna().mean()*100)`

```
funder                     0.0
gps_height                 0.0
installer                  0.0
longitude                  0.0
latitude                   0.0
basin                      0.0
region                     0.0
district_code              0.0
lga                        0.0
ward                       0.0
population                 0.0
public_meeting             0.0
permit                     0.0
construction_year          0.0
extraction_type_group      0.0
management                 0.0
payment                    0.0
water_quality              0.0
quantity                   0.0
source                     0.0
waterpoint_type            0.0
status_group               0.0
decade                     0.0
installer_classified       0.0
funder_classified          0.0
dtype: float64
```

**population**

[272]: ```
#Getting the mean and median
train_data1.loc[train_data1['population']!=0].describe()
```

[272]:
|       | gps_height    | longitude     | latitude      | district_code | population    | \ |
|-------|---------------|---------------|---------------|---------------|---------------|---|
| count | 38019.000000  | 38019.000000  | 38019.000000  | 38019.000000  | 38019.000000  |   |
| mean  | 969.889634    | 36.074387     | -6.139781     | 6.299456      | 281.087167    |   |
| std   | 612.544787    | 2.586779      | 2.737733      | 11.303334     | 564.687660    |   |

```
min       -90.000000    29.607122    -11.649440     1.000000      1.000000
25%       347.000000    34.715340     -8.388839     2.000000     40.000000
50%      1135.000000    36.706815     -5.750877     3.000000    150.000000
75%      1465.000000    37.940149     -3.597016     5.000000    324.000000
max      2770.000000    40.345193     -1.042375    67.000000  30500.000000


       construction_year
count        38019.000000
mean          1961.399721
std            263.994165
min              0.000000
25%           1986.000000
50%           2000.000000
75%           2008.000000
max           2013.000000
```

[273]: 
```python
#Replacing the population that is 0 with mean
train_data1['population'].replace(to_replace = 0 , value =281, inplace=True)
```

[274]: 
```python
print(train_data1.isna().mean()*100)
```

```
funder                  0.0
gps_height              0.0
installer               0.0
longitude               0.0
latitude                0.0
basin                   0.0
region                  0.0
district_code           0.0
lga                     0.0
ward                    0.0
population              0.0
public_meeting          0.0
permit                  0.0
construction_year       0.0
extraction_type_group   0.0
management              0.0
payment                 0.0
water_quality           0.0
quantity                0.0
source                  0.0
waterpoint_type         0.0
status_group            0.0
decade                  0.0
installer_classified    0.0
funder_classified       0.0
dtype: float64
```

### 0.4.11 Data without nullvalues for EDA analysis

## 0.5 Exploratory Data Analysis

### 0.5.1 Univariate Analysis

**Categorical Variables**

```
[275]: # Define the categorical columns for analysis
       categorical_cols = ['basin', 'region', 'public_meeting', 'permit',
        ↪'extraction_type_group', 'management', 'payment', 'water_quality',
        ↪'quantity', 'source', 'waterpoint_type', 'decade', 'installer_classified',
        ↪'funder_classified']

       # Define the number of columns in each row of subplots
       num_cols_per_row = 2

       # Calculate the number of rows needed for subplots
       num_rows = (len(categorical_cols) + num_cols_per_row - 1) // num_cols_per_row

       # Create subplots with adjusted width and height
       fig, axs = plt.subplots(num_rows, num_cols_per_row, figsize=(20, num_rows * 7))

       # Flatten the axis array for easier iteration
       axs = axs.flatten()

       # Loop through each categorical column and create count plots
       for i, col in enumerate(categorical_cols):
           sns.countplot(x=col, data=train_data1, ax=axs[i])
           axs[i].set_title(f'Count Plot of {col}')
           axs[i].tick_params(axis='x', rotation=45)

       # Hide any extra subplots if the number of columns is not a multiple of
        ↪num_cols_per_row
       for j in range(len(categorical_cols), num_rows * num_cols_per_row):
           fig.delaxes(axs[j])

       plt.tight_layout()
       plt.show()
```

**Numeric Variables**

```
[276]:  # Assuming train_data1 is your DataFrame
        population_array = np.array(train_data1['population'])

        # Set the size of the plot
        plt.figure(figsize=(40, 15))

        # Plot the boxplot
        plt.subplot(2, 2, 2)
        sns.boxplot(y=population_array, color='green')
        plt.title('Boxplot of Population ')
        plt.ylabel('Population')

        # Show the plot
        plt.show()


        gps_height_array = np.array(train_data1['gps_height'])
        # Plot the boxplot

        # Set the size of the plot
        plt.figure(figsize=(40, 15))
        plt.subplot(2, 2, 2)
        sns.boxplot(y=gps_height_array , color='green')
        plt.title('Boxplot of gps_height ')
        plt.ylabel('gps_height')

        # Show the plot
        plt.show()
```

Boxplot of gps_height

### 0.5.2 Bivariate Analysis

**Categorical variables**

```
[278]:  # Define the categorical and numeric columns for analysis
        categorical_cols = ['basin', 'region', 'public_meeting', 'permit',
         ↪'extraction_type_group', 'management', 'payment', 'water_quality',
         ↪'quantity', 'source', 'waterpoint_type', 'decade', 'installer_classified',
         ↪'funder_classified', 'district_code']
        numeric_cols = ['gps_height', 'population']


        # Define the number of columns in each row of subplots
        num_cols_per_row = 1


        # Calculate the number of rows needed for subplots
        num_rows_cat = (len(categorical_cols) + num_cols_per_row - 1) //
         ↪num_cols_per_row
        num_rows_num = (len(numeric_cols) + num_cols_per_row - 1) // num_cols_per_row


        # Inspect the unique values in the status_group column
        unique_status_groups = train_data1['status_group'].unique()
        print(f"Unique status groups: {unique_status_groups}")


        # Create a color palette for the status groups
        # Ensure these keys match the unique values exactly
        palette = {
            'functional': 'green',
            'non functional': 'red',
            'functional needs repair': 'orange'
        }


        # Create subplots for categorical vs. target variable
```

```
fig, axs_cat = plt.subplots(num_rows_cat, num_cols_per_row, figsize=(15,␣
 ↪num_rows_cat * 6))
axs_cat = axs_cat.flatten()

for i, col in enumerate(categorical_cols):
    if col in ['installer_classified', 'funder_classified']:
        # Filter out 'others' values
        filtered_data = train_data1[train_data1[col] != 'Others']
    else:
        filtered_data = train_data1

    sns.countplot(x=col, hue='status_group', data=filtered_data, ax=axs_cat[i],␣
 ↪palette=palette)
    axs_cat[i].set_title(f'Count Plot of {col} by Status Group')
    axs_cat[i].tick_params(axis='x', rotation=45)
    axs_cat[i].legend(title='Status Group', loc='upper right')

# Hide any extra subplots if the number of columns is not a multiple of␣
 ↪num_cols_per_row
for j in range(len(categorical_cols), num_rows_cat * num_cols_per_row):
    fig.delaxes(axs_cat[j])

plt.tight_layout()
plt.show()
```

Unique status groups: ['functional' 'non functional' 'functional needs repair']

Count Plot of basin by Status Group


Count Plot of region by Status Group


Count Plot of public_meeting by Status Group


Count Plot of permit by Status Group


Count Plot of extraction_type_group by Status Group


Count Plot of management by Status Group


Count Plot of payment by Status Group


Count Plot of water_quality by Status Group


Count Plot of quantity by Status Group


Count Plot of source by Status Group


Count Plot of waterpoint_type by Status Group


Count Plot of decade by Status Group


Count Plot of installer_classified by Status Group


Count Plot of funder_classified by Status Group

Count Plot of ... by Status Group

**Numerical variables**

```
[279]:  # Define a color mapping for the status groups
        status_colors = {
            'functional': 'green',
            'non functional': 'red',
            'functional needs repair': 'orange'
        }

        # Map status_group to colors
        train_data1['color'] = train_data1['status_group'].map(status_colors)

        # Aggregate population by status group (sum)
        population_sum = train_data1.groupby('status_group')['population'].sum().
          ↪reset_index()


        # Ensure the order of bars matches the status_colors keys order
        status_order = ['functional', 'non functional', 'functional needs repair']
        colors = [status_colors[status] for status in status_order]

        # Create a bar plot
        plt.figure(figsize=(20, 8))  # Adjust figure size for better spacing
        sns.barplot(x='status_group', y='population', data=population_sum,␣
          ↪palette=colors, order=status_order)
        plt.title('Total Population by Status Group')
        plt.xlabel('Status Group')
        plt.ylabel('Total Population')
        plt.xticks(rotation=45)

        # Show the plot
        plt.show()

        # Filter out rows where population is zero or missing, if needed
        train_data1 = train_data1[train_data1['population'] > 0]

        # Aggregate population by status group (mean)
        population_mean = train_data1.groupby('status_group')['population'].mean().
          ↪reset_index()

        # Create a bar plot
        plt.figure(figsize=(20, 8))  # Adjust figure size for better spacing
        sns.barplot(x='status_group', y='population',␣
          ↪data=population_mean,palette=colors)
        plt.title('Average Population by Status Group')
```

```python
plt.xlabel('Status Group')
plt.ylabel('Average Population')
plt.xticks(rotation=45)

# Show the plot
plt.show()


# Ensure 'gps_height' and 'status_group' are in the dataframe
if 'gps_height' in train_data1.columns and 'status_group' in train_data1.
 ↪columns:
    # Filter out rows where gps_height is missing or zero, if needed
    train_data1 = train_data1[train_data1['gps_height'] != 0]

    # Aggregate gps_height by status group (mean)
    gps_height_mean = train_data1.groupby('status_group')['gps_height'].mean().
 ↪reset_index()

    # Create a bar plot
    plt.figure(figsize=(20, 8))  # Adjust figure size for better spacing
    sns.barplot(x='status_group', y='gps_height', data=gps_height_mean,␣
 ↪palette=colors)
    plt.title('Average GPS Height by Status Group')
    plt.xlabel('Status Group')
    plt.ylabel('Average GPS Height')
    plt.xticks(rotation=45)

    # Show the plot
    plt.show()
else:
    print("Columns 'gps_height' or 'status_group' not found in the dataset.")

    # Define a color mapping for the status groups
status_colors = {
    'functional': 'green',
    'non functional': 'red',
    'functional needs repair': 'orange'
}

# Map status_group to colors
train_data1['color'] = train_data1['status_group'].map(status_colors)

# Create a scatter plot
fig, ax = plt.subplots(figsize=(20, 15))  # Increase the figure size

# Plot the scatter plot
```

```
scatter = ax.scatter(train_data1['longitude'], train_data1['latitude'],
 ↪c=train_data1['color'], alpha=0.6, s=10)

# Add basemap from OpenStreetMap
# Note: This step assumes you have contextily installed and can fetch the
 ↪basemap.
# Uncomment the below lines if contextily is available and installed
#import contextily as ctx
#ctx.add_basemap(ax, crs='EPSG:4326', source=ctx.providers.OpenStreetMap.Mapnik)

# Set title and labels
plt.title('Distribution of Water Point Status Group in Tanzania', fontsize=20)
plt.xlabel('Longitude', fontsize=15)
plt.ylabel('Latitude', fontsize=15)

# Create a custom legend
legend_elements = [
    Line2D([0], [0], marker='o', color='w', label='Functional',
 ↪markerfacecolor='green', markersize=15),
    Line2D([0], [0], marker='o', color='w', label='Non Functional',
 ↪markerfacecolor='red', markersize=15),
    Line2D([0], [0], marker='o', color='w', label='Functional Needs Repair',
 ↪markerfacecolor='orange', markersize=15)
]

ax.legend(handles=legend_elements, loc='upper right', fontsize=15)

plt.show()
```



49

Average Population by Status Group



Average GPS Height by Status Group

Distribution of Water Point Status Group in Tanzania

### 0.5.3 Multivariate Analysis

```
[280]:  # Select only numeric columns for correlation computation
        numeric_cols = train_data1.select_dtypes(include=['float64', 'int64']).columns
        numeric_data = train_data1[numeric_cols]

        # Compute the correlation matrix
        correlation_matrix = numeric_data.corr()

        # Plot the heatmap
        plt.figure(figsize=(12, 8))
        sns.heatmap(correlation_matrix, annot=True, fmt='.2f', cmap='coolwarm')
        plt.title('Correlation Matrix')
        plt.show()
```

Correlation Matrix

**Observations**

1. Government-funded wells often exhibit a higher likelihood of being non-functional, highlighting a need for improved oversight or maintenance practices in these projects.

2. Areas with higher populations tend to have a greater number of functional wells, indicating a correlation between population density and well functionality.

3. Certain areas show a higher probability of accessing clean water, particularly those situated near good water basins, highlighting the importance of geographical location in water quality.

4. Despite being one of the most densely populated cities, Dar es Salaam has a significant portion (35%) of clean water sources classified as non-functional, indicating challenges in maintaining water infrastructure.

5. Iringa, an important area, has a notable number of non-functional water points with soft water, suggesting potential issues with water quality or infrastructure maintenance in this region.

6. Water points installed by central government and district councils also show a tendency towards non-functionality, indicating potential systemic issues in water infrastructure management at the governmental level.

7. While gravity-based extraction is the most common type, hand pumps, which are less efficient, rank second. This suggests a need for authorities to focus on upgrading or maintaining pumping infrastructure, particularly for gravity-based systems that are naturally reliant on gravitational forces.

8. Some water points with sufficient and soft water are non-functional, indicating that water quality alone may not guarantee well functionality and that other factors like maintenance play a crucial role.

9. Recent years have seen a higher proportion of functional wells compared to older ones, but there are still functional wells that require repair. This underscores the importance of timely maintenance to prevent functional wells from deteriorating into non-functional ones.

10. Many water wells with ample water resources are non-functional, highlighting potential issues with infrastructure or operational aspects rather than water availability.

**Recommendations**

1. Targeted Maintenance: Prioritize maintenance in densely populated areas and near good water basins, focusing on government-funded wells and central installations.

2. Water Quality Focus: Improve water quality monitoring and treatment, especially in areas with soft water but high non-functionality rates like Iringa.

3. Pumping Infrastructure Investment: Upgrade pumping infrastructure, particularly hand pumps and gravity systems, to enhance efficiency and reduce non-functional wells.

## 0.6 Modelling

[281]: `train_data1.info()`

```
<class 'pandas.core.frame.DataFrame'>
Index: 38962 entries, 0 to 59399
Data columns (total 26 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   funder                38962 non-null  object
 1   gps_height            38962 non-null  int64
 2   installer             38962 non-null  object
 3   longitude             38962 non-null  float64
 4   latitude              38962 non-null  float64
 5   basin                 38962 non-null  object
 6   region                38962 non-null  object
 7   district_code         38962 non-null  int64
 8   lga                   38962 non-null  object
 9   ward                  38962 non-null  object
 10  population            38962 non-null  int64
 11  public_meeting        38962 non-null  bool
 12  permit                38962 non-null  bool
 13  construction_year     38962 non-null  int64
 14  extraction_type_group 38962 non-null  object
```

```
15  management          38962 non-null  object
16  payment             38962 non-null  object
17  water_quality       38962 non-null  object
18  quantity            38962 non-null  object
19  source              38962 non-null  object
20  waterpoint_type     38962 non-null  object
21  status_group        38962 non-null  object
22  decade              38962 non-null  object
23  installer_classified 38962 non-null object
24  funder_classified   38962 non-null  object
25  color               38962 non-null  object
dtypes: bool(2), float64(2), int64(4), object(18)
memory usage: 7.5+ MB
```

### 0.6.1 Drop columns that are not necessary for our modeling

[282]:
```
train_data1.
 ↪drop(columns=['funder','installer','construction_year','color','lga','ward'],inplace=True
 ↪)
```

[283]:
```
train_data1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 38962 entries, 0 to 59399
Data columns (total 20 columns):
 #   Column               Non-Null Count  Dtype
---  ------               --------------  -----
 0   gps_height           38962 non-null  int64
 1   longitude            38962 non-null  float64
 2   latitude             38962 non-null  float64
 3   basin                38962 non-null  object
 4   region               38962 non-null  object
 5   district_code        38962 non-null  int64
 6   population           38962 non-null  int64
 7   public_meeting       38962 non-null  bool
 8   permit               38962 non-null  bool
 9   extraction_type_group 38962 non-null object
 10  management           38962 non-null  object
 11  payment              38962 non-null  object
 12  water_quality        38962 non-null  object
 13  quantity             38962 non-null  object
 14  source               38962 non-null  object
 15  waterpoint_type      38962 non-null  object
 16  status_group         38962 non-null  object
 17  decade               38962 non-null  object
 18  installer_classified 38962 non-null  object
 19  funder_classified    38962 non-null  object
dtypes: bool(2), float64(2), int64(3), object(13)
```

```
memory usage: 5.7+ MB
```

```
[284]:  # Convert public_meeting, permit to 0 and 1
        # Convert 'permit' and 'public_meeting' to binary (0 and 1)
        train_data1['permit'] = train_data1['permit'].map({True: 1, False: 0})
        train_data1['public_meeting'] = train_data1['public_meeting'].map({True: 1,␣
        ↪False: 0})
        train_data1.head()
```

```
[284]:      gps_height  longitude   latitude                    basin    region  \
        0         1390  34.938093  -9.856322              Lake Nyasa    Iringa
        1         1399  34.698766  -2.147466            Lake Victoria      Mara
        2          686  37.460664  -3.821329                  Pangani   Manyara
        3          263  38.486161 -11.155298  Ruvuma / Southern Coast    Mtwara
        10          62  39.209518  -7.034139              Wami / Ruvu     Pwani

            district_code  population  public_meeting  permit extraction_type_group  \
        0               5         109               1       0               gravity
        1               2         280               1       1               gravity
        2               4         250               1       1               gravity
        3              63          58               1       1            submersible
        10             43         345               1       0            submersible

                  management         payment water_quality     quantity  \
        0                vwc     pay annually          soft       enough
        1                wug        never pay          soft  insufficient
        2                vwc   pay per bucket          soft       enough
        3                vwc        never pay          soft          dry
        10  private operator        never pay         salty       enough

                        source              waterpoint_type    status_group decade  \
        0               spring          communal standpipe       functional    90s
        1   rainwater harvesting         communal standpipe       functional    10s
        2                  dam  communal standpipe multiple      functional    00s
        3          machine dbh  communal standpipe multiple  non functional    80s
        10         machine dbh                        other      functional    10s

           installer_classified funder_classified
        0                Others            Others
        1                Others            Others
        2          world vision            Others
        3                Others            Unicef
        10               Others            Others
```

```
[285]:  train_data1['status_group'].value_counts()
```

```
[285]: status_group
       functional                21790
       non functional            14618
       functional needs repair    2554
       Name: count, dtype: int64
```

### 0.6.2 Convert target to ternary values

```
[286]: target_status_group = {'functional':0,
                               'non functional': 2,
                               'functional needs repair': 1}
       train_data1['status_group'] = train_data1['status_group'].
        ↪replace(target_status_group)
```

```
[287]: train_data1['status_group'].value_counts()
```

```
[287]: status_group
       0    21790
       2    14618
       1     2554
       Name: count, dtype: int64
```

### 0.6.3 Having my numerical,target and my categorical columns

```
[288]: categorical_columns =␣
        ↪['basin','region','extraction_type_group','management','payment','water_quality','quantity'
                  ␣
        ↪'source','waterpoint_type','decade','installer_classified','funder_classified']
```

```
[289]: numerical_columns =␣
        ↪['gps_height','longitude','latitude','district_code','population','public_meeting','permit']
```

```
[290]: target='status_group'
```

### 0.6.4 Logistics Regression Model

**Assumptions**

1. Linearity: Assumes a linear relationship between the independent variables and the log-odds of the dependent variable.
2. Independence of Errors: Assumes that the errors (residuals) of the observations are independent.
3. No Multicollinearity: Assumes that the independent variables are not highly correlated with each other. Multicollinearity can inflate the variance of coefficient estimates and make the model unstable.
4. Homoscedasticity: Assumes constant variance of errors.

5. Binary Outcome: Assumes the outcome variable is binary (can be extended to multinomial logistic regression for multiple classes).

```
[291]: from sklearn.model_selection import train_test_split, RandomizedSearchCV,
        ↪GridSearchCV
       from sklearn.preprocessing import StandardScaler, OneHotEncoder
       from sklearn.compose import ColumnTransformer
       from sklearn.pipeline import Pipeline
       from sklearn.linear_model import LogisticRegression
       from sklearn.metrics import classification_report, accuracy_score,
        ↪confusion_matrix,balanced_accuracy_score
       from sklearn.impute import SimpleImputer
       from scipy.stats import uniform
       from sklearn.svm import SVC
```

```
[292]: # Separate features and target
       X = train_data1[categorical_columns + numerical_columns]
       y = train_data1[target]
```

```
[293]: # Define preprocessor
       preprocessor = ColumnTransformer(
           transformers=[
               ('num', StandardScaler(), numerical_columns),
               ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_columns)
           ])

       # Define pipeline
       pipeline = Pipeline(steps=[
           ('preprocessor', preprocessor),
           ('classifier', LogisticRegression(max_iter=1000))
       ])

       # Split the data into train and test sets
       X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
        ↪random_state=42)

       # Fit the model
       pipeline.fit(X_train, y_train)

       # Predict on the test set
       y_pred = pipeline.predict(X_test)

       y_pred_train = pipeline.predict(X_train)

       # Evaluate the model
       print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
       print('Test Accuracy:', accuracy_score(y_test, y_pred))
```

```python
print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))


# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=pipeline.
  ↪classes_, yticklabels=pipeline.classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```
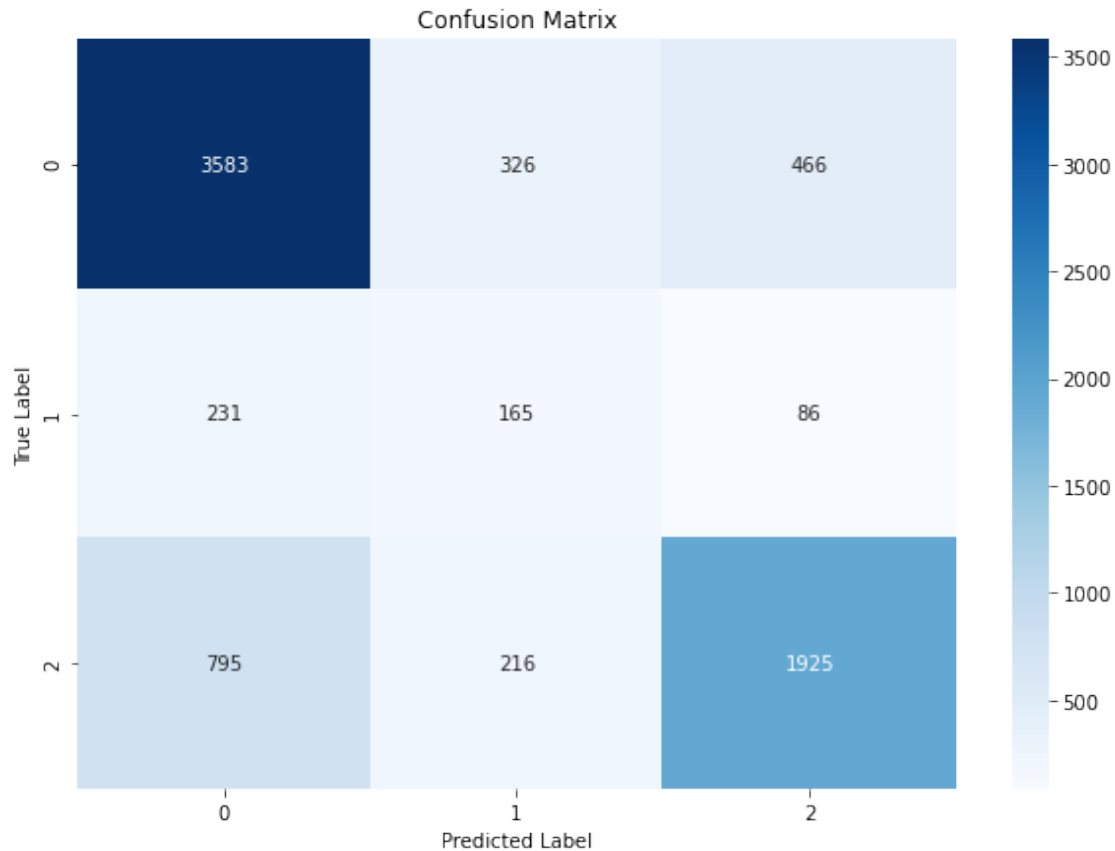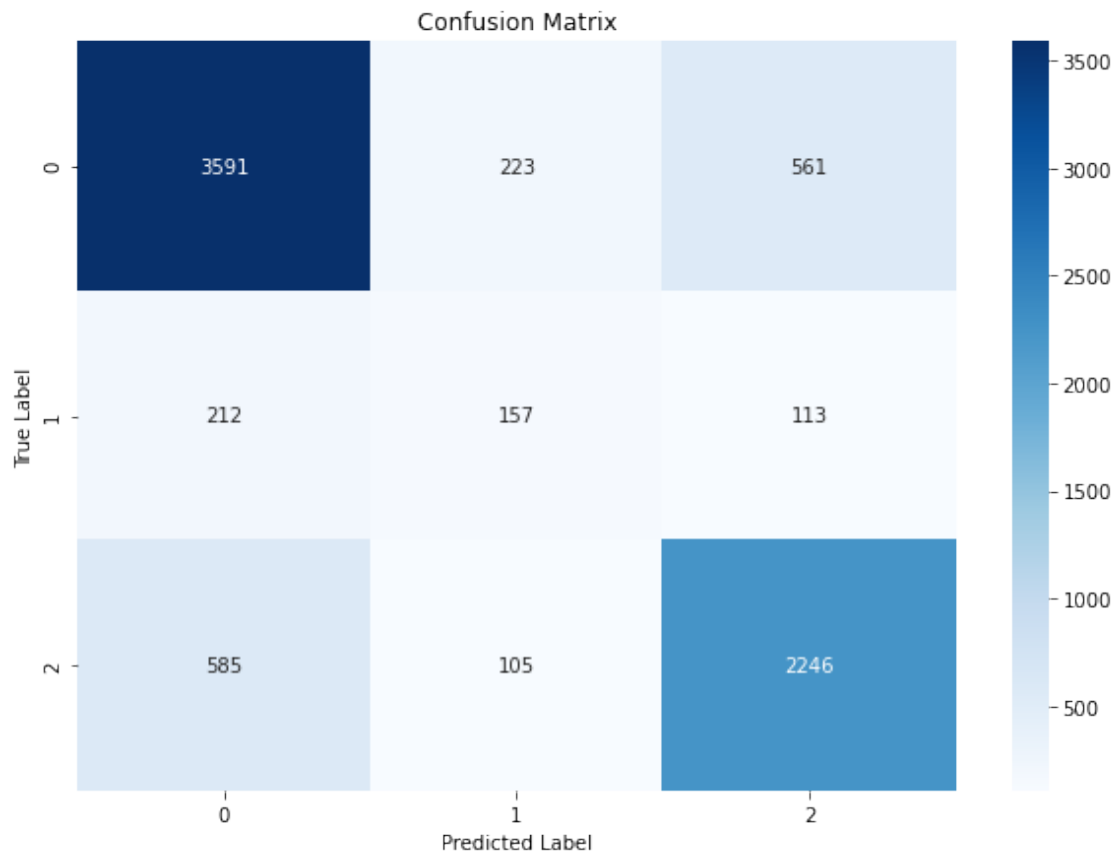
Train Accuracy: 0.7608842118771857
Test Accuracy: 0.7633773899653535
Balance Train Accuracy: 0.5641972764919815
Balance Test Accuracy: 0.5654770868968072

### 0.6.5 Tuned Logistic Regression Model

```python
# Define preprocessor
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numerical_columns),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_columns)
    ])

# Define pipeline
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier',
 ↪LogisticRegression(class_weight='balanced',solver='liblinear',max_iter=1000))
])

# Split the data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=42)

# Manually tune parameters
best_score = 0
best_params = {}
for C in [0.001, 0.01, 0.1, 1, 10]:
    for penalty in ['l1', 'l2']:
        try:
            # Set parameters
            pipeline.set_params(classifier__C=C, classifier__penalty=penalty)
            # Train model
            pipeline.fit(X_train, y_train)
            # Evaluate model
            score = accuracy_score(y_test, pipeline.predict(X_test))

            # Update best parameters
            if score > best_score:
                best_score = score
                best_params = {'C': C, 'penalty': penalty}
        except ValueError as e:
            pass


# Train the best model
pipeline.set_params(classifier__C=best_params['C'],
 ↪classifier__penalty=best_params['penalty'])
pipeline.fit(X_train, y_train)

# Predict and evaluate
```

[294]:

```python
y_pred = pipeline.predict(X_test)
y_pred_train = pipeline.predict(X_train)

# Evaluate the model
print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
print('Test Accuracy:', accuracy_score(y_test, y_pred))
print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))


# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=pipeline.
 ↪classes_, yticklabels=pipeline.classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```

```
Train Accuracy: 0.7338701915364625
Test Accuracy: 0.727960990632619
Balance Train Accuracy: 0.6191026542009768
Balance Test Accuracy: 0.6056496769924631
```

Confusion Matrix

### 0.6.6 Decision tree Model

**Assumptions**

1. Independence of Observations: Assumes that the observations in the dataset are independent of each other.
2. No Assumption on Feature Distribution: Decision Trees do not make assumptions about the distribution of the data.
3. Sufficient Data: Requires a large enough dataset to adequately split and create meaningful branches.
4. Minimal Preprocessing: Can handle both numerical and categorical data, and does not require data scaling or normalization.
5. Non-Linearity: Can capture non-linear relationships between features and the target variable.

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.tree import DecisionTreeClassifier
```

```python
from sklearn.metrics import classification_report,
 ↪confusion_matrix,balanced_accuracy_score
import seaborn as sns
import matplotlib.pyplot as plt
```

[296]:
```python
# Preprocess the categorical and numerical columns
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numerical_columns),
        ('cat', OneHotEncoder(), categorical_columns)
    ])

# Create the decision tree pipeline
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', DecisionTreeClassifier(random_state=42))
])

# Train the model
pipeline.fit(X_train, y_train)

# Predict on the test set
y_pred = pipeline.predict(X_test)
y_pred_train = pipeline.predict(X_train)

# Evaluate the model
print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
print('Test Accuracy:', accuracy_score(y_test, y_pred))
print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))


# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
 ↪xticklabels=pipeline.classes_, yticklabels=pipeline.classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```
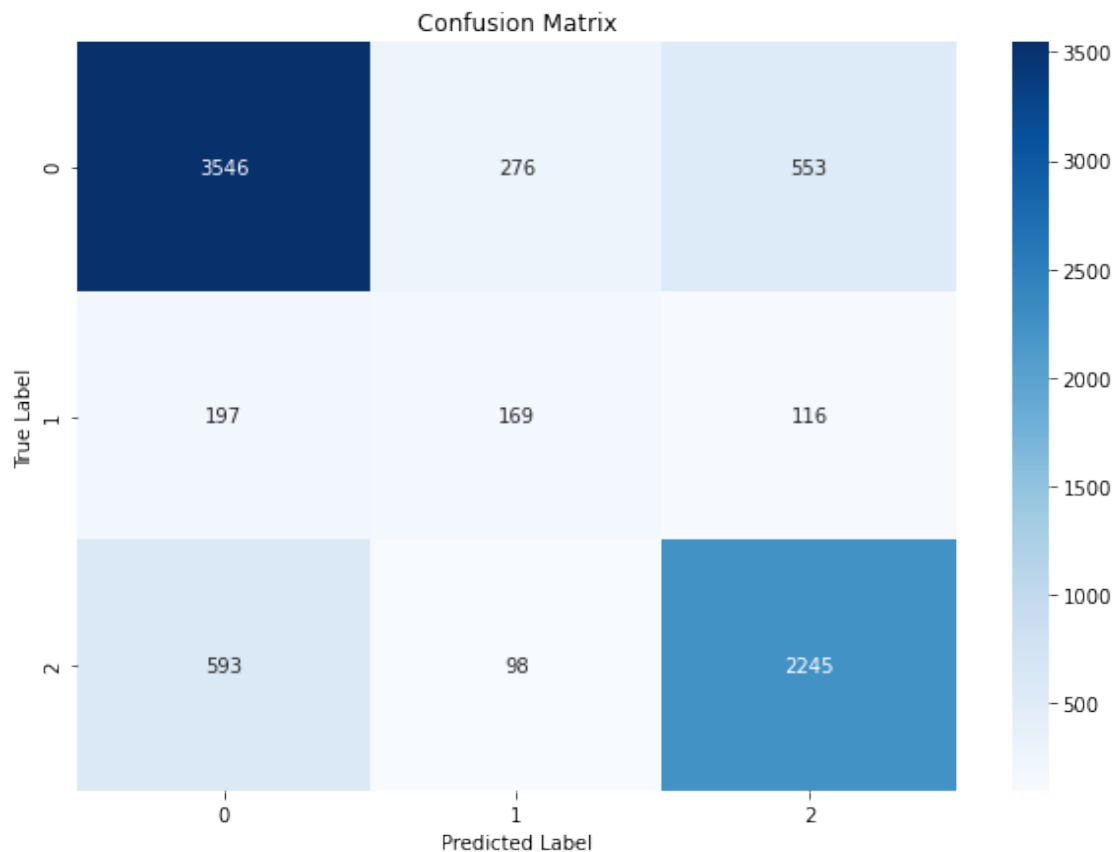
```
Train Accuracy: 1.0
Test Accuracy: 0.7691518029000385
Balance Train Accuracy: 1.0
```

Balance Test Accuracy: 0.6371708390335455



## 0.7 Tuned Decision Tree Model

```
[297]:  # Preprocess the categorical and numerical columns
        preprocessor = ColumnTransformer(
            transformers=[
                ('num', StandardScaler(), numerical_columns),
                ('cat', OneHotEncoder(), categorical_columns)
            ])

        # Create the decision tree pipeline
        pipeline = Pipeline(steps=[
            ('preprocessor', preprocessor),
            ('classifier',
          →DecisionTreeClassifier(class_weight='balanced',random_state=42))
        ])

        # Define a range of hyperparameters to search over
        criteria = ['gini', 'entropy']
```

```python
max_depths = [None, 10, 20, 30]
min_samples_splits = [2, 5, 10]
min_samples_leaves = [1, 2, 4]

best_score = 0
best_params = {}

# Manually tune hyperparameters
for criterion in criteria:
    for max_depth in max_depths:
        for min_samples_split in min_samples_splits:
            for min_samples_leaf in min_samples_leaves:
                # Set parameters
                pipeline.set_params(
                    classifier__criterion=criterion,
                    classifier__max_depth=max_depth,
                    classifier__min_samples_split=min_samples_split,
                    classifier__min_samples_leaf=min_samples_leaf
                )
                # Train model
                pipeline.fit(X_train, y_train)
                # Evaluate model
                score = accuracy_score(y_test, pipeline.predict(X_test))
                # Update best parameters
                if score > best_score:
                    best_score = score
                    best_params = {
                        'criterion': criterion,
                        'max_depth': max_depth,
                        'min_samples_split': min_samples_split,
                        'min_samples_leaf': min_samples_leaf
                    }

# Train the best model
pipeline.set_params(
    classifier__criterion=best_params['criterion'],
    classifier__max_depth=best_params['max_depth'],
    classifier__min_samples_split=best_params['min_samples_split'],
    classifier__min_samples_leaf=best_params['min_samples_leaf']
)
pipeline.fit(X_train, y_train)

# Predict and evaluate
y_pred = pipeline.predict(X_test)
y_pred_train = pipeline.predict(X_train)

# Evaluate the model
```

```python
print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
print('Test Accuracy:', accuracy_score(y_test, y_pred))
print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))


# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',␣
 ↪xticklabels=pipeline.classes_, yticklabels=pipeline.classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```
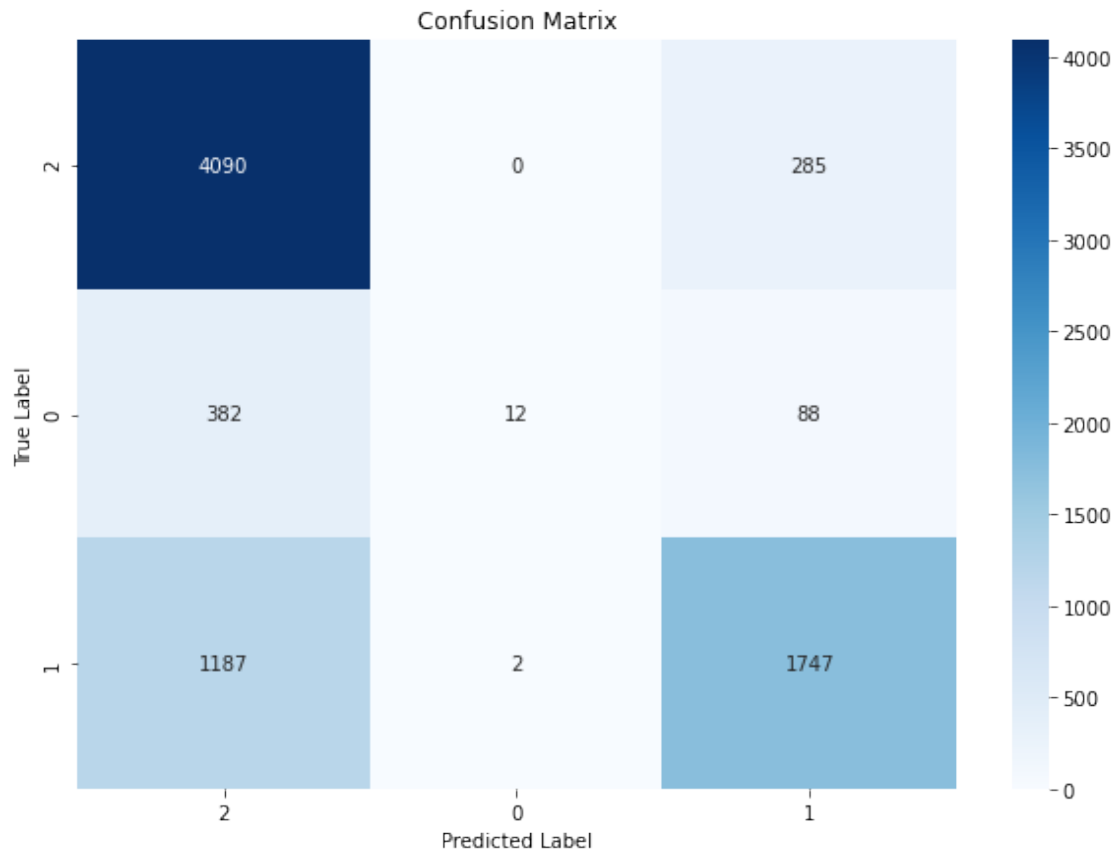
Train Accuracy: 0.9904392184542333
Test Accuracy: 0.7647889131271655
Balance Train Accuracy: 0.9933755576465688
Balance Test Accuracy: 0.6419274896400157

### 0.7.1 Support Vector Machine Model

**Assumptions**

1. Linearly Separable Data (for Linear SVM): Assumes that the data is linearly separable if using a linear kernel.
2. Margin Maximization: SVM tries to find the hyperplane that maximizes the margin between the classes.
3. Kernel Trick (for Non-Linear SVM): Assumes that the kernel function can transform the data into a higher-dimensional space where it is linearly separable.
4. Feature Scaling: Assumes that the data is scaled properly. SVMs are sensitive to the scale of the input features.
5. Independence of Observations: Assumes that the observations are independent of each other.

```python
[298]: # Preprocess the categorical and numerical columns
       preprocessor = ColumnTransformer(
           transformers=[
               ('num', StandardScaler(), numerical_columns),
               ('cat', OneHotEncoder(), categorical_columns)
           ])

       # Create the SVM pipeline
       pipeline = Pipeline(steps=[
           ('preprocessor', preprocessor),
           ('classifier', SVC(kernel='linear', random_state=42))
       ])

       # Train the model
       pipeline.fit(X_train, y_train)

       # Predict on the test set
       y_pred = pipeline.predict(X_test)
       y_pred_train = pipeline.predict(X_train)


       # Evaluate the model
       print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
       print('Test Accuracy:', accuracy_score(y_test, y_pred))
       print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
       print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))


       # Confusion Matrix
       conf_matrix = confusion_matrix(y_test, y_pred)
```

```python
# Get the unique classes from the target
classes = y_train.unique()

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',␣
  ↪xticklabels=classes, yticklabels=classes)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```

Train Accuracy: 0.7499117713112388
Test Accuracy: 0.7505453612216091
Balance Train Accuracy: 0.5193383759702246
Balance Test Accuracy: 0.5182602187912374



## 0.7.2   Random Forest Machine Learning Model

**Assumptions**

67

1. Independence of Observations: Assumes that the observations in the dataset are independent of each other.
2. No Assumption on Distribution: Unlike some models, Random Forest does not make strong assumptions about the distribution of the data.
3. Feature Importance: Assumes that some features are more important than others, and the model will try to identify and leverage these important features.
4. Sufficient Data: Requires a sufficiently large dataset to build diverse and effective trees.
5. Minimal Preprocessing: Can handle missing values and does not require much preprocessing (like scaling) of the data.

```python
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import StandardScaler, OneHotEncoder,RobustScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline,make_pipeline
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, balanced_accuracy_score,
 ↪confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns
import category_encoders as ce

# Assuming X and y are your features and target variables
X = pd.DataFrame(X, columns=numerical_columns + categorical_columns)
y = pd.Series(y, name=target)

# Split the data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=42)

# choosing scaler and encoder
scaler=RobustScaler()
encoder = ce.TargetEncoder(cols=categorical_columns)

# putting numeric columns to scaler and categorical to encoder
num_transformer = make_pipeline(scaler)
cat_transformer = make_pipeline(encoder)

# Preprocess the categorical and numerical columns
preprocessor = ColumnTransformer(
     transformers=[('num', num_transformer, numerical_columns),
                   ('cat', cat_transformer, categorical_columns)])

# Create the Random Forest pipeline
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', RandomForestClassifier(random_state=42))
```

```python
])

# Train the model
pipeline.fit(X_train, y_train)

# Predict on the test set
y_pred = pipeline.predict(X_test)
y_pred_train = pipeline.predict(X_train)

# Evaluate the model
print('Train Accuracy:', accuracy_score(y_train, y_pred_train))
print('Test Accuracy:', accuracy_score(y_test, y_pred))
print('Balance Train Accuracy:', balanced_accuracy_score(y_train, y_pred_train))
print('Balance Test Accuracy:', balanced_accuracy_score(y_test, y_pred))

# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
  ↪xticklabels=pipeline.classes_, yticklabels=pipeline.classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```

```
Train Accuracy: 0.9999679168404505
Test Accuracy: 0.8179135121262672
Balance Train Accuracy: 0.9998391248391248
Balance Test Accuracy: 0.6558696180171392
```

Confusion Matrix

**Tuning Random forest Model- My Final Model**

```
[300]: import pandas as pd
       import seaborn as sns
       import matplotlib.pyplot as plt
       from sklearn.compose import ColumnTransformer
       from sklearn.ensemble import RandomForestClassifier
       from sklearn.metrics import accuracy_score, balanced_accuracy_score,
        ↪confusion_matrix
       from sklearn.model_selection import train_test_split
       from sklearn.pipeline import make_pipeline
       from sklearn.preprocessing import RobustScaler
       import category_encoders as ce


       # Split the data into training and testing sets
       X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
        ↪random_state=42)

       # Identify numerical and categorical columns
```

```python
numerical_columns = X.select_dtypes(include=['int64', 'float64']).columns.
  ↪tolist()
categorical_columns = X.select_dtypes(include=['object']).columns.tolist()

# Choosing scaler and encoder
scaler = RobustScaler()
encoder = ce.TargetEncoder(cols=categorical_columns)

# Putting numeric columns to scaler and categorical to encoder
num_transformer = make_pipeline(scaler)
cat_transformer = make_pipeline(encoder)

# Getting together our scaler and encoder with preprocessor
preprocessor = ColumnTransformer(
    transformers=[('num', num_transformer, numerical_columns),
                  ('cat', cat_transformer, categorical_columns)]
)

# Set RandomForestClassifier with initial parameters
rf = RandomForestClassifier(n_estimators=100, random_state=42, n_jobs=-1,
                            criterion='entropy', max_features='sqrt',
                            min_samples_split=10, class_weight='balanced')

# Giving all values to pipeline
pipeline = make_pipeline(preprocessor, rf)

# Define hyperparameter ranges
n_estimators_options = [100, 200]
max_depth_options = [None, 10, 20]
min_samples_split_options = [2, 5]
min_samples_leaf_options = [2, 4]
max_features_options = ['sqrt']
bootstrap_options = [True, False]

best_score = 0
best_params = {}

# Loop through all combinations of hyperparameters
for n_estimators in n_estimators_options:
    for max_depth in max_depth_options:
        for min_samples_split in min_samples_split_options:
            for min_samples_leaf in min_samples_leaf_options:
                for max_features in max_features_options:
                    for bootstrap in bootstrap_options:
                        # Update the model in the pipeline
                        pipeline.
  ↪set_params(randomforestclassifier__n_estimators=n_estimators,
```

```python
      ↪randomforestclassifier__max_depth=max_depth,
      ↪
      ↪randomforestclassifier__min_samples_split=min_samples_split,
      ↪
      ↪randomforestclassifier__min_samples_leaf=min_samples_leaf,
      ↪
      ↪randomforestclassifier__max_features=max_features,
      ↪
      ↪randomforestclassifier__bootstrap=bootstrap)

                        # Train the model
                        pipeline.fit(X_train, y_train)

                        # Predict on the test set
                        y_test_pred = pipeline.predict(X_test)

                        # Evaluate the model
                        score = balanced_accuracy_score(y_test, y_test_pred)

                        # If the current score is better than the best score,␣
      ↪update best score and best params
                        if score > best_score:
                            best_score = score
                            best_params = {
                                'n_estimators': n_estimators,
                                'max_depth': max_depth,
                                'min_samples_split': min_samples_split,
                                'min_samples_leaf': min_samples_leaf,
                                'max_features': max_features,
                                'bootstrap': bootstrap
                            }

# Print the best hyperparameters
print("Best Hyperparameters:", best_params)

# Train the final model with the best hyperparameters on the entire training set
pipeline.
 ↪set_params(randomforestclassifier__n_estimators=best_params['n_estimators'],
                randomforestclassifier__max_depth=best_params['max_depth'],
      ↪
 ↪randomforestclassifier__min_samples_split=best_params['min_samples_split'],
      ↪
 ↪randomforestclassifier__min_samples_leaf=best_params['min_samples_leaf'],
      ↪
 ↪randomforestclassifier__max_features=best_params['max_features'],
```

```python
                        randomforestclassifier__bootstrap=best_params['bootstrap'])
pipeline.fit(X_train, y_train)

# Predictions on train set
y_pred = pipeline.predict(X_train)

# Predictions on test set
y_pred_test = pipeline.predict(X_test)

# Evaluate the model
train_accuracy = accuracy_score(y_train, y_pred)
test_accuracy = accuracy_score(y_test, y_pred_test)
balance_train_accuracy = balanced_accuracy_score(y_train, y_pred)
balance_test_accuracy = balanced_accuracy_score(y_test, y_pred_test)
print(f"Train Accuracy: {train_accuracy:.4f}")
print(f"Test Accuracy: {test_accuracy:.4f}")
print(f"Balanced Train Accuracy: {balance_train_accuracy:.4f}")
print(f"Balanced Test Accuracy: {balance_test_accuracy:.4f}")

# Confusion Matrix
cm = confusion_matrix(y_test, y_pred_test)

# Plot the confusion matrix
plt.figure(figsize=(10, 7))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=pipeline.
 ↪named_steps['randomforestclassifier'].classes_, yticklabels=pipeline.
 ↪named_steps['randomforestclassifier'].classes_)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```
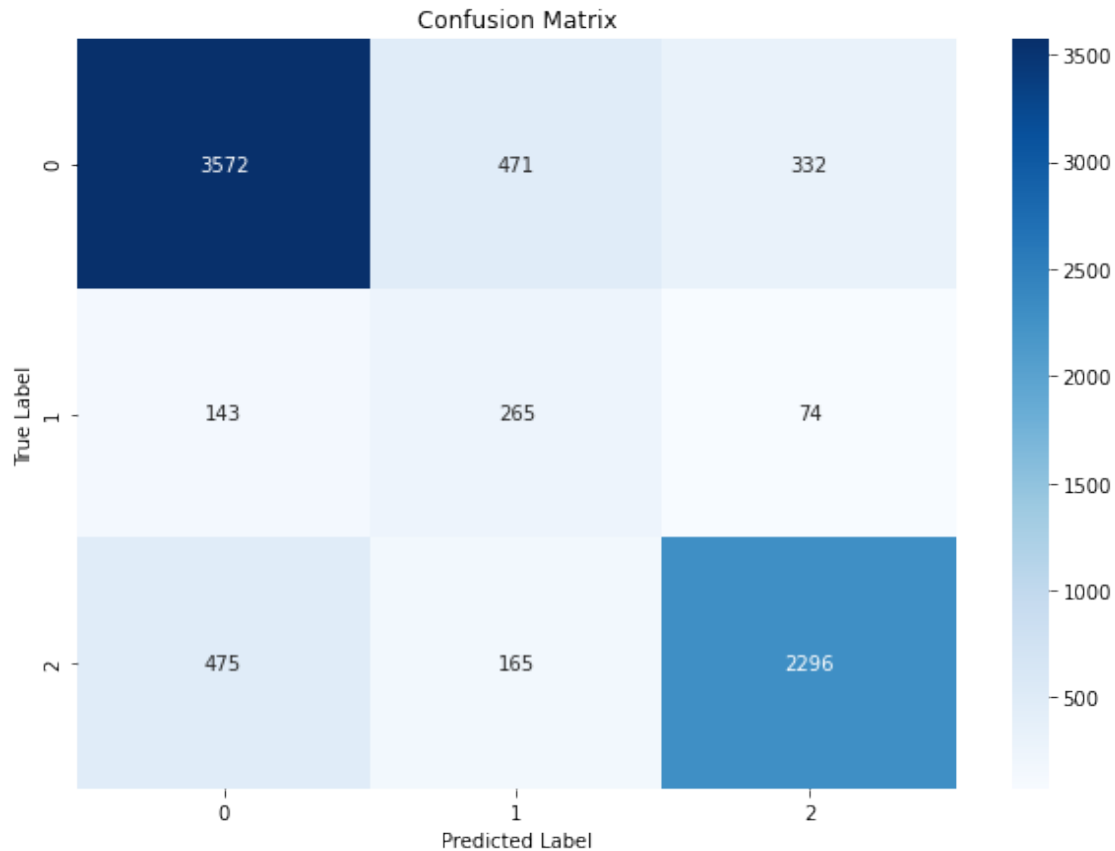
```
Best Hyperparameters: {'n_estimators': 200, 'max_depth': 20,
'min_samples_split': 2, 'min_samples_leaf': 4, 'max_features': 'sqrt',
'bootstrap': False}
Train Accuracy: 0.9025
Test Accuracy: 0.7870
Balanced Train Accuracy: 0.9310
Balanced Test Accuracy: 0.7161
```

Confusion Matrix

### 0.7.3 Conclusion on my final model- Random Forest

The final model is a RandomForestClassifier with tuned hyperparameters trained on a dataset split into training and testing sets. The hyperparameters were tuned using an exhaustive search through various combinations to find the ones that maximize the balanced accuracy score on the test set. The best hyperparameters identified were as follows: 'n_estimators': 200, 'max_depth': 20, 'min_samples_split': 2, 'min_samples_leaf': 4, 'max_features': 'sqrt', 'bootstrap': False.

The model achieved a decent level of performance, with a test accuracy of 78.70% and a balanced test accuracy of 71.61%. These metrics indicate that the model generalizes reasonably well to unseen data and is not overfitting excessively to the training set. The balanced accuracy score is particularly useful in scenarios where classes are imbalanced, as it takes into account the imbalance and provides a more reliable measure of overall model performance.

The confusion matrix plot visualizes how well the model is predicting each class. It shows the number of true positives, true negatives, false positives, and false negatives for each class, allowing for a deeper understanding of the model's strengths and weaknesses in classification. Overall, the final model appears to be a solid choice for the given dataset and task.

### 0.7.4 Predictive Analysis of Tanzanian Water Well Conditions

**Model Performance** The classifier built to predict the condition of water wells in Tanzania achieved a test accuracy of 78.70% and a balanced test accuracy of 71.61%. These metrics suggest that the model performs reasonably well in identifying the condition of water wells based on features

**Important Features** The most important features identified by the model include the type of pump used, the installation year, and possibly other geographic or environmental factors. Understanding these key features can help stakeholders prioritize maintenance and repair efforts for water wells.

**Useful Predictions** For an NGO focused on locating wells needing repair, the model's predictions can be highly valuable. By identifying non-functional or deteriorating wells accurately, the NGO can allocate resources more efficiently and effectively, ensuring that clean water access is maintained or restored where needed most.

**Recommendations for Stakeholders:**

1. Modify Input Variables: Based on the model's insights, stakeholders could consider modifying certain input variables. For example, investing in newer pump technologies or improving maintenance schedules for wells installed in specific years could lead to better overall well conditions.

2. Target Results: The model can help stakeholders set specific targets for well conditions. By analyzing patterns in non-functional wells, they can influence how new wells are built, ensuring they are more resilient and require less frequent repairs.

3. Geographical Considerations: Considering geographic or environmental factors that influence well conditions can further enhance the model's predictive capabilities. For instance, areas with certain soil types or rainfall patterns may require different pump types or maintenance strategies.

In conclusion, the predictive model offers valuable insights into the condition of Tanzanian water wells, aiding stakeholders in making informed decisions regarding maintenance, repair, and future well construction strategies.

### 0.7.5 Next Steps

1. Validation and Deployment of Model: Validate the predictive model using additional datasets or real-time data to ensure its accuracy and reliability. Once validated, deploy the model for ongoing monitoring and prediction of water well conditions.

2. Actionable Insights Implementation: Implement actionable insights derived from the EDA analysis, such as prioritizing maintenance in high-population areas, improving water quality monitoring, and investing in pumping infrastructure. Collaborate with stakeholders and authorities to translate these insights into practical initiatives.

3. Continuous Improvement: Continuously evaluate and improve the model based on feedback and new data. Incorporate feedback from field teams, stakeholders, and ongoing data collection to refine the model's predictive capabilities and enhance decision-making related to water well management