

Detection of Non-functional Water Wells Using Machine Learning Algorithms



Overview

Tanzania, as a developing country, struggles with providing clean water to its population of over 57,000,000. There are many wells (water points) already established in the country, but some are in need of repair while others have failed altogether. The Government of Tanzania is looking to find patterns in non-functional wells to influence how new wells are built.

Business Problem

Tanzania faces challenges in ensuring access to clean water for its population due to non-functional water wells. We will be looking to predict patterns in non-functional wells to inform more robust construction methods for new wells.

Data Understanding

Dataset

Driven Data - Tanzanian Water Wells

- [Labels \(<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/>\)](https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/)
- [Values \(<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/>\)](https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/)

Target

- functional : the well is operational and there are no repairs needed
- functional needs repair : the well is operational, but needs repairs
- non functional : the well is not operational

Features

- amount_tsh : Total static head (amount water available to well)
- date_recorded : The date the row was entered
- funder : Who funded the well
- gps_height : Altitude of the well
- installer : Organization that installed the well
- longitude : GPS coordinate
- latitude : GPS coordinate
- wpt_name : Name of the well if there is one
- num_private : Private use or not
- basin : Geographic water basin
- subvillage : Geographic location
- region : Geographic location
- region_code : Geographic location (coded)
- district_code : Geographic location (coded)
- lga : Geographic location
- ward : Geographic location
- population : Population around the well
- public_meeting : True/False
- recorded_by : Group entering this row of data
- scheme_management : Who operates the well
- scheme_name : Who operates the well
- permit : If the well is permitted
- construction_year : Year the well was constructed
- extraction_type : The kind of extraction the well uses
- extraction_type_group : The kind of extraction the well uses
- extraction_type_class : The kind of extraction the well uses
- management : How the well is managed
- management_group : How the well is managed
- payment : What the water costs
- payment_type : What the water costs
- water_quality : The quality of the water
- quality_group : The quality of the water
- quantity : The quantity of water
- quantity_group : The quantity of water
- source : The source of the water
- source_type : The source of the water
- source_class : The source of the water
- waterpoint_type : The kind of well
- waterpoint_type_group : The kind of well

Limitations

- Do not know what types of repairs are needed
- Not able to use the well age, due to missing construction years
- Without full population information, we do not know the supply needs

Data Preparation

Import and Read Datasets

```
In [1]: 1 # Import standard packages
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 %matplotlib inline
7 import datetime
8
9 # Model Selection
10 from sklearn.model_selection import train_test_split, GridSearchCV,
11
12 # Classification Models
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.tree import DecisionTreeClassifier
15 from sklearn.dummy import DummyClassifier
16 from sklearn.ensemble import RandomForestClassifier
17
18 from sklearn.metrics import plot_confusion_matrix, accuracy_score,
19 from sklearn.metrics import classification_report, confusion_matrix
20 from sklearn.metrics import roc_curve, auc, roc_auc_score
21
22 # Scalers
23 from sklearn.impute import SimpleImputer
24 from sklearn.preprocessing import StandardScaler
25
26 # Categorical Create Dummies
27 from sklearn.preprocessing import OneHotEncoder
28
29 # Column Transformer
30 from sklearn.compose import ColumnTransformer
31
32 # Pipeline
33 from sklearn.pipeline import Pipeline
34
35 # Base
36 from sklearn.base import BaseEstimator, TransformerMixin
37
38 import warnings
39 warnings.filterwarnings('ignore')
```

```
In [2]: 1 # Load the datasets
2 features_data = 'data/WellWaterData.csv'
3 target_data = 'data/TargetData.csv'
4
5 features = pd.read_csv(features_data)
6 target = pd.read_csv(target_data)
7
```

Merging Datasets

```
In [3]: 1 # Checking for unique IDs in both datasets to ensure they match
2 unique_ids_features = features['id'].nunique()
3 unique_ids_target = target['id'].nunique()
4
5 unique_ids_features, unique_ids_target
6
```

Out[3]: (59400, 59400)

```
In [4]: 1 # Merging the datasets on the 'id' column
2 merged_data = pd.merge(features, target, on='id')
3
4 # Displaying the first few rows of the merged dataset
5 merged_data.head()
6
```

Out[4]:

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latitude	wl
0	69572	6000.0	2011-03-14	Roman	1390	Roman	34.938093	-9.856322	
1	8776	0.0	2013-03-06	Grumeti	1399	GRUMETI	34.698766	-2.147466	
2	34310	25.0	2013-02-25	Lottery Club	686	World vision	37.460664	-3.821329	
3	67743	0.0	2013-01-28	Unicef	263	UNICEF	38.486161	-11.155298	Near
4	19728	0.0	2011-07-13	Action In A	0	Artisan	31.130847	-1.825359	

5 rows × 41 columns

In [5]: 1 merged_data.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 59400 entries, 0 to 59399
Data columns (total 41 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           55765 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         59400 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 55523 non-null    object  
 21  scheme_name       31234 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  
 40  status_group        59400 non-null    object  
dtypes: float64(3), int64(7), object(31)
memory usage: 19.0+ MB
```

Clean

Preparing the merged dataset for feature exploration and how they relate to the 'status_group' target variable.

```
In [6]: 1 # Converting 'date_recorded' to datetime
2 merged_data['date_recorded'] = pd.to_datetime(merged_data['date_rec
```

Addressing Categorical Features with Parent and Subgroup Columns

```
In [7]: 1 grouped = merged_data.groupby(['extraction_type_class', 'extraction_type_group'])
2 grouped
```

```
Out[7]: extraction_type_class    extraction_type_group    extraction_type
gravity                           gravity                  gravity
26780
handpump                         afridev                afridev
1770
2400
98
8154
ga      2
85
229
48
3670
motorpump                         mono                   mono
2865
90
32
other
6430
rope pump                         rope pump              other - rope pump
451
submersible                        submersible            ksb
1415
4764
wind-powered                       wind-powered           windmill
117
dtype: int64
```

```
In [8]: 1 grouped = merged_data.groupby(['management_group', 'management']).size()
```

```
In [9]: 1 grouped = merged_data.groupby(['waterpoint_type_group', 'waterpoint_type']).size()
```

Handling Missing Values

```
In [10]: 1 # Calculating the percentage of zero values for each column
2 zero_value_percentages = {}
3 for column in merged_data.columns:
4     zero_count = (merged_data[column] == 0).sum()
5     zero_value_percentages[column] = (zero_count / len(merged_data))
6
7
8 zero_value_percentages
```

```
Out[10]: {'id': 0.0016835016835016834,
'amount_tsh': 70.09932659932659,
'date_recorded': 0.0,
'funder': 0.0,
'gps_height': 34.40740740740741,
'installer': 0.0,
'longitude': 3.05050505050505,
'latitude': 0.0,
'wpt_name': 0.0,
'num_private': 98.72558922558923,
'basin': 0.0,
'subvillage': 0.0,
'region': 0.0,
'region_code': 0.0,
'district_code': 0.038720538720538725,
'lga': 0.0,
'ward': 0.0,
'population': 35.994949494949495,
'public_meeting': 8.51010101010101,
'recorded_by': 0.0,
'scheme_management': 0.0,
'scheme_name': 0.0,
'permit': 29.44781144781145,
'construction_year': 34.86363636363636,
'extraction_type': 0.0,
'extraction_type_group': 0.0,
'extraction_type_class': 0.0,
'management': 0.0,
'management_group': 0.0,
'payment': 0.0,
'payment_type': 0.0,
'water_quality': 0.0,
'quality_group': 0.0,
'quantity': 0.0,
'quantity_group': 0.0,
'source': 0.0,
'source_type': 0.0,
'source_class': 0.0,
'waterpoint_type': 0.0,
'waterpoint_type_group': 0.0,
'status_group': 0.0}
```

```
In [11]: 1 # Calculating the percentage of missing values in each column
2 missing_values = merged_data.isnull().mean() * 100
3 missing_values = missing_values[missing_values > 0].sort_values(ascending=True)
4
5 missing_values
6
7
```

```
Out[11]: scheme_name      47.417508
scheme_management    6.526936
installer           6.153199
funder              6.119529
public_meeting       5.612795
permit               5.144781
subvillage          0.624579
dtype: float64
```

Since 'scheme_management', 'installer', 'funder', and 'permit' have less than 7% missing values and are potentially relevant, replacing them is a good option.

```
In [12]: 1 # Replacing the missing values
2 for column in ['scheme_management', 'installer', 'funder', 'permit']:
3     merged_data[column].fillna('Unknown', inplace=True)
4
5 # Sanity check on missing values
6 remaining_missing_values = merged_data.isnull().sum()
7 remaining_missing_values[remaining_missing_values > 0]
```

```
Out[12]: subvillage      371
public_meeting    3334
scheme_name       28166
dtype: int64
```

Dropping Columns

```
In [13]: 1 # Dropping additional group columns
2 # Dropping irrelevant columns
3 # Dropping features with missing values over 45%
4 features_dropped = merged_data.drop(columns=['quantity_group', 'extraction_type',
5                                              'waterpoint_type_group',
6                                              'payment', 'water_quality',
7                                              'region_code', 'public_meeting',
8                                              'wpt_name', 'district_code',
9                                              'ward', 'subvillage'],
10                                             errors='ignore')
```

Removing Duplicates

In [14]: 1 features_dropped.duplicated().sum()

Out[14]: 685

In [15]: 1 data_dedup = features_dropped.drop_duplicates()
2
3 # Rechecking for duplicates
4 new_duplicate_count = data_dedup.duplicated().sum()
5 new_duplicate_count

Out[15]: 0

Feature Engineering

In [16]: 1 # Extracting year and month from 'date_recorded'
2 data_dedup['year_recorded'] = data_dedup['date_recorded'].dt.year
3 data_dedup['month_recorded'] = data_dedup['date_recorded'].dt.month
4
5 data_dedup[['longitude', 'gps_height', 'construction_year', 'year_recorded', 'month_recorded']]

Out[16]:

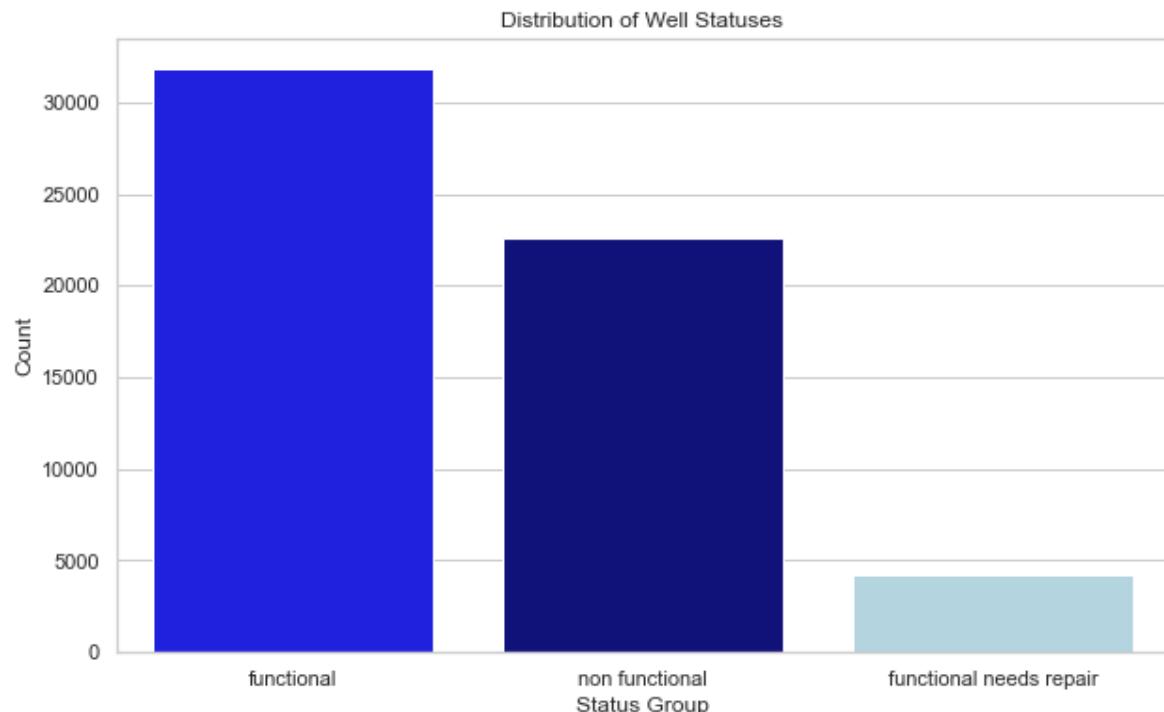
	longitude	gps_height	construction_year	year_recorded	month_recorded
0	34.938093	1390	1999	2011	3
1	34.698766	1399	2010	2013	3
2	37.460664	686	2009	2013	2
3	38.486161	263	1986	2013	1
4	31.130847	0	0	2011	7

In [17]: 1 # Adding 'well_age' feature
2 data_dedup['well_age'] = data_dedup.apply(lambda row: 0 if row['construction_year'] < 1900 else row['year_recorded'] - row['construction_year'])
3
4
5

Creating Binary Target Column

First, we'll need to understand the distribution of the target variable and visualize the proportion of functional vs. non-functional wells.

```
In [18]: 1 sns.set(style="whitegrid")
2
3 # Defining palette colors
4 palette_colors = {'functional': 'blue', 'non functional': 'darkblue'
5
6 # Plotting the distribution of well statuses
7 plt.figure(figsize=(10, 6))
8 sns.countplot(x='status_group', palette=palette_colors, data=data_d)
9 plt.title('Distribution of Well Statuses')
10 plt.ylabel('Count')
11 plt.xlabel('Status Group')
12 plt.show()
```



To address differing opinions on how to condense our target into a binary column, we will create two separate binary target columns and assess the better performer on our baseline model.

status_binary:

- Class 0 = non-functional & functional needs repair
- Class 1 = functional

status_binary_reversed:

- Class 0 = non-functional
- Class 1 = functional & functional needs repair

```
In [19]: 1 # Binary encoding of the 'status_group' column
2
3 # 'functional' is assigned 1 and 'non functional' or 'functional ne
4 data_dedup['status_binary'] = data_dedup['status_group'].apply(lambda
5
6 # 'non-functional' is assigned 0 and 'functional' or 'functional ne
7 data_dedup['status_binary_reversed'] = data_dedup['status_group'].a
8
```

Creating Master Dataset

```
In [20]: 1 # Creating master dataset with all values of 'well_age' greater tha
2 # to eliminate negative values where 'recorded_year' was likely lis
3 master_data = data_dedup[data_dedup['well_age'] >= 0]
4 master_data.info()
5
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 58706 entries, 0 to 59399
Data columns (total 26 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   amount_tsh        58706 non-null   float64
 1   date_recorded    58706 non-null   datetime64[ns]
 2   funder            58706 non-null   object 
 3   gps_height       58706 non-null   int64  
 4   installer         58706 non-null   object 
 5   longitude         58706 non-null   float64
 6   latitude          58706 non-null   float64
 7   basin             58706 non-null   object 
 8   region            58706 non-null   object 
 9   population        58706 non-null   int64  
 10  scheme_management 58706 non-null   object 
 11  permit             58706 non-null   object 
 12  construction_year 58706 non-null   int64  
 13  extraction_type_class 58706 non-null   object 
 14  management         58706 non-null   object 
 15  payment_type       58706 non-null   object 
 16  quality_group      58706 non-null   object 
 17  quantity            58706 non-null   object 
 18  source_type         58706 non-null   object 
 19  waterpoint_type    58706 non-null   object 
 20  status_group        58706 non-null   object 
 21  year_recorded      58706 non-null   int64  
 22  month_recorded     58706 non-null   int64  
 23  well_age            58706 non-null   int64  
 24  status_binary       58706 non-null   int64  
 25  status_binary_reversed 58706 non-null   int64  
dtypes: datetime64[ns](1), float64(3), int64(8), object(14)
memory usage: 12.1+ MB
```

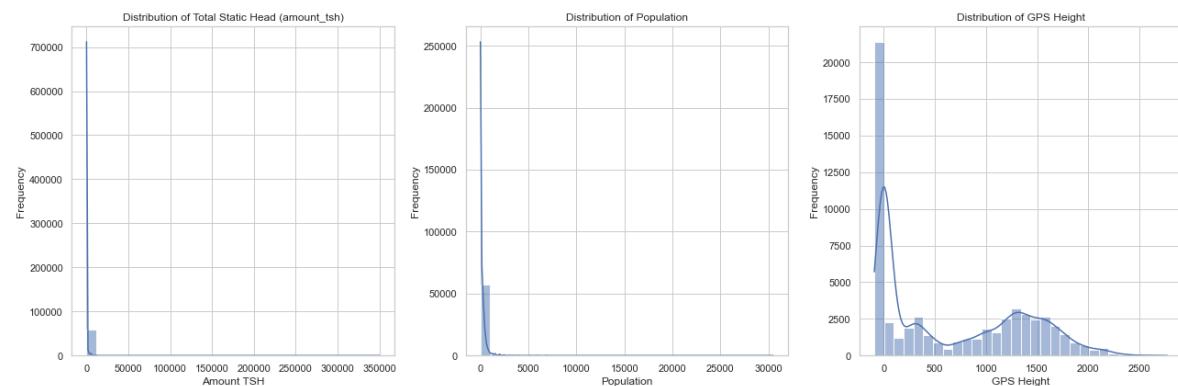
```
In [21]: 1 # Saving master dataset to csv
          2 master_data.to_csv('data/master_data.csv', index=False)
```

EDA

Data Distribution Visualizations

Using histograms, we'll plot the distribution of key numeric variables like amount_tsh, population, and gps_height.

```
In [22]: 1 sns.set_style("whitegrid")
          2
          3 # Creating histograms for 'amount_tsh', 'population', and 'gps_height'
          4 fig, axes = plt.subplots(1, 3, figsize=(18, 6))
          5
          6 # Plot for 'amount_tsh'
          7 sns.histplot(master_data['amount_tsh'], bins=30, ax=axes[0], kde=True)
          8 axes[0].set_title('Distribution of Total Static Head (amount_tsh)')
          9 axes[0].set_xlabel('Amount TSH')
         10 axes[0].set_ylabel('Frequency')
         11
         12 # Plot for 'population'
         13 sns.histplot(master_data['population'], bins=30, ax=axes[1], kde=True)
         14 axes[1].set_title('Distribution of Population')
         15 axes[1].set_xlabel('Population')
         16 axes[1].set_ylabel('Frequency')
         17
         18 # Plot for 'gps_height'
         19 sns.histplot(master_data['gps_height'], bins=30, ax=axes[2], kde=True)
         20 axes[2].set_title('Distribution of GPS Height')
         21 axes[2].set_xlabel('GPS Height')
         22 axes[2].set_ylabel('Frequency')
         23
         24 plt.tight_layout()
         25 plt.show()
         26
```



Observations

These distributions reflect the high number of zero values in our dataset for these features.

Total Static Head (amount_tsh): The distribution appears to be highly skewed to the right, indicating that most wells have a low static head value.

Population: This distribution is also right-skewed, showing that most wells serve a relatively small population, with fewer points serving larger populations.

GPS Height: The distribution is more varied, indicating a range of elevations at which wells are located.

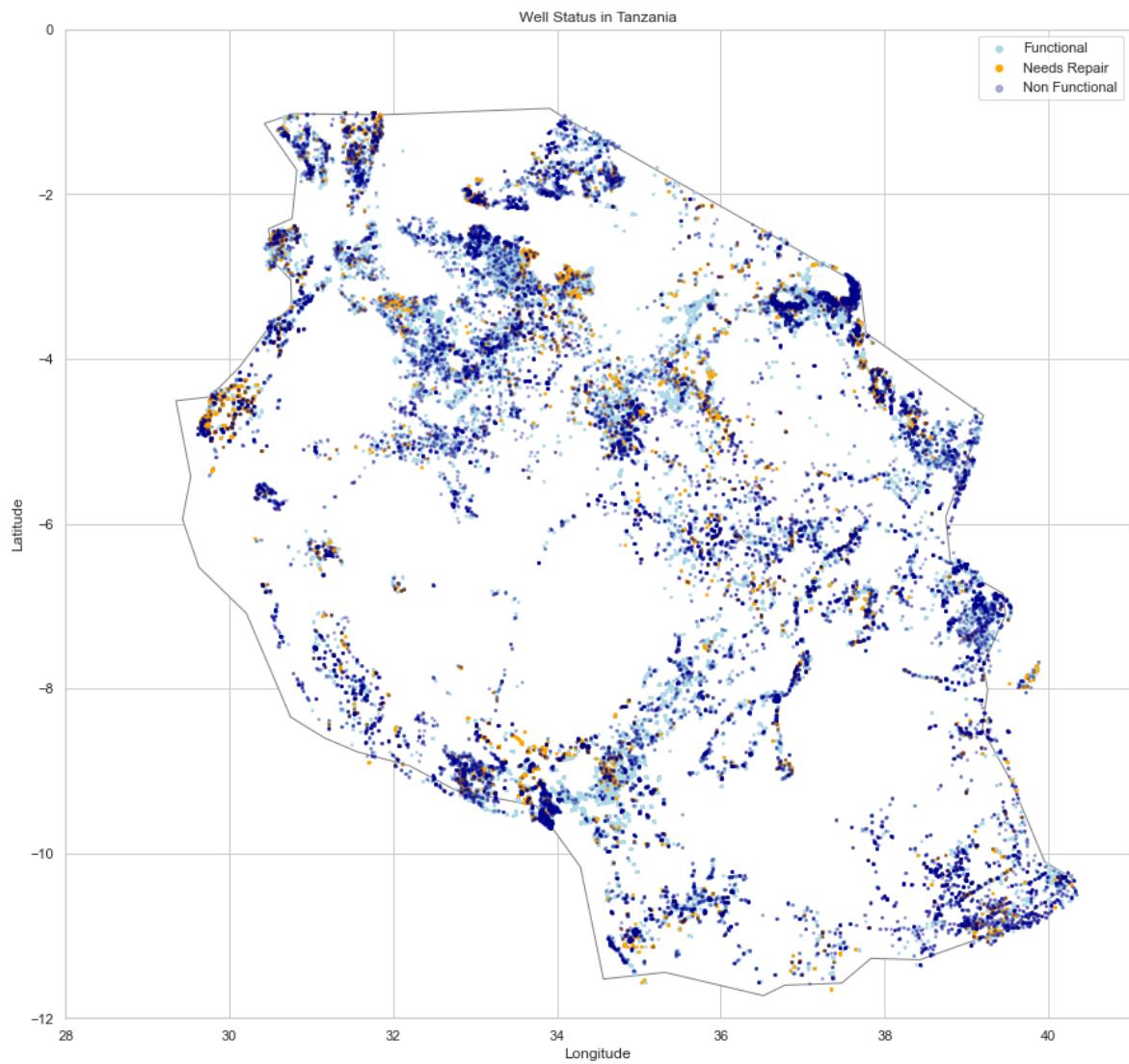
Geographical Analysis

We'll create a geographical plot using latitude and longitude to see if there is any geographical pattern in the status of wells.

```
In [23]: 1 !pip install geopandas
```

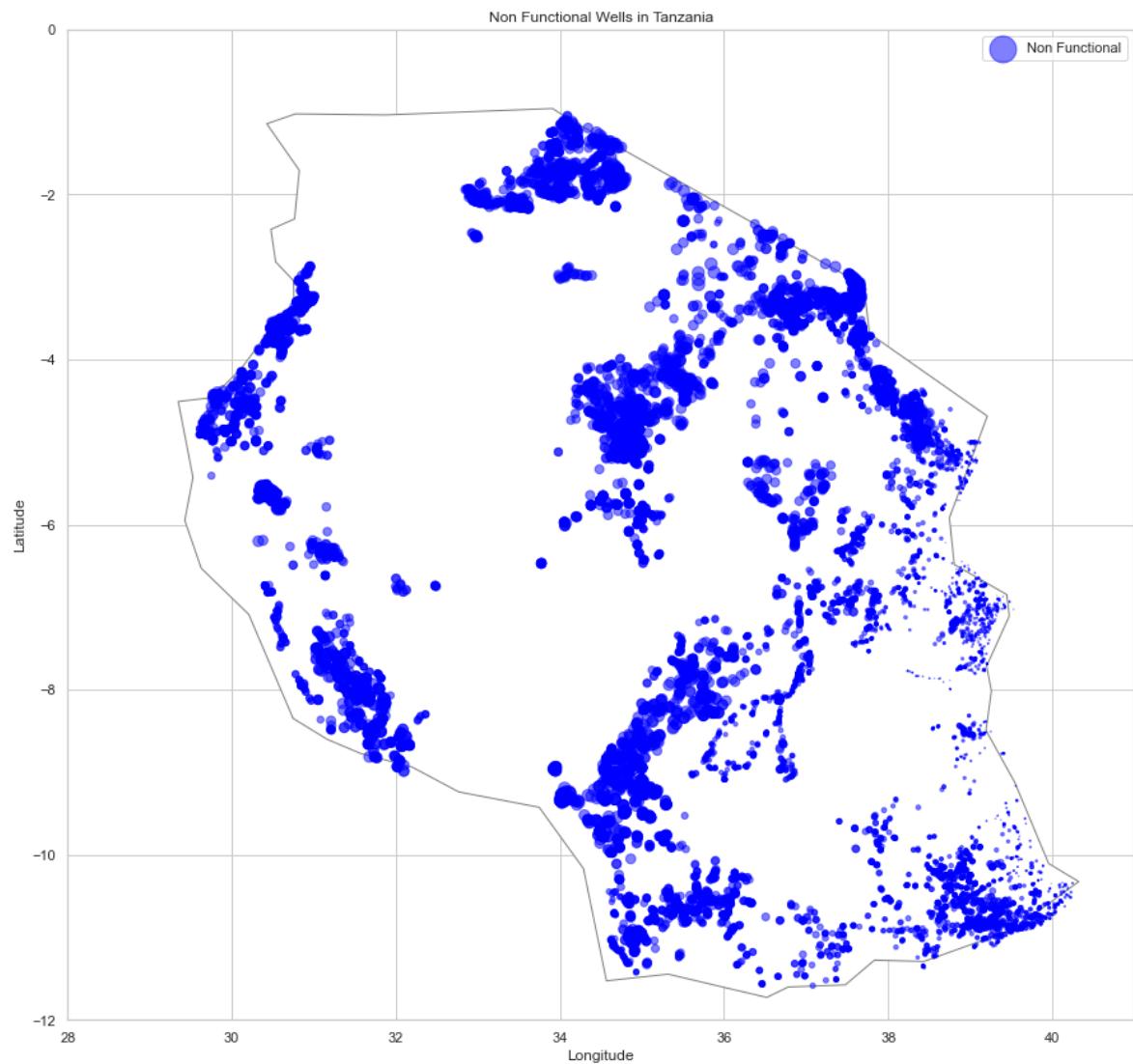
```
Requirement already satisfied: geopandas in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (0.13.2)
Requirement already satisfied: fiona>=1.8.19 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from geopandas) (1.9.5)
Requirement already satisfied: packaging in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from geopandas) (20.4)
Requirement already satisfied: pandas>=1.1.0 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from geopandas) (1.1.3)
Requirement already satisfied: shapely>=1.7.1 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from geopandas) (2.0.2)
Requirement already satisfied: pyproj>=3.0.1 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from geopandas) (3.5.0)
Requirement already satisfied: click~=8.0 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (8.1.7)
Requirement already satisfied: cligj>=0.5 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (0.7.2)
Requirement already satisfied: importlib-metadata; python_version < "3.10" in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (2.0.0)
Requirement already satisfied: six in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (1.15.0)
Requirement already satisfied: click-plugins>=1.0 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (1.1.1)
Requirement already satisfied: attrs>=19.2.0 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (20.2.0)
Requirement already satisfied: certifi in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (2020.6.20)
Requirement already satisfied: setuptools in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from fiona>=1.8.19->geopandas) (50.3.0.post20201103)
Requirement already satisfied: pyparsing>=2.0.2 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from packaging->geopandas) (2.4.7)
Requirement already satisfied: python-dateutil>=2.7.3 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from pandas>=1.1.0->geopandas) (2.8.1)
Requirement already satisfied: pytz>=2017.2 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from pandas>=1.1.0->geopandas) (2020.1)
Requirement already satisfied: numpy>=1.15.4 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from pandas>=1.1.0->geopandas) (1.18.5)
Requirement already satisfied: zipp>=0.5 in /Users/goknurkaya/anaconda3/envs/learn-env/lib/python3.8/site-packages (from importlib-metadata; python_version < "3.10"->fiona>=1.8.19->geopandas) (3.3.0)
```

```
In [24]: 1 import geopandas
2 import matplotlib.pyplot as plt
3
4 # Create GeoDataFrame
5 gdf = geopandas.GeoDataFrame(
6     master_data, geometry=geopandas.points_from_xy(master_data.longitude,
7
8 # Assigning 'status_group'
9 functional = gdf[gdf['status_group'] == 'functional']
10 repair = gdf[gdf['status_group'] == 'functional needs repair']
11 non_functional = gdf[gdf['status_group'] == 'non functional']
12
13 # Load world shapefile
14 world_shapefile_path = 'data/ne_110m_admin_0_countries/ne_110m_adm1.shp'
15 world = geopandas.read_file(world_shapefile_path)
16
17 # Filter for Tanzania
18 fig, ax = plt.subplots(figsize=(15, 15))
19 base = world[world.ADMIN == 'United Republic of Tanzania'].plot(color='white')
20
21 # Scatter plots for each category
22 ax.scatter(functional['longitude'], functional['latitude'], c='lightblue', alpha=0.5)
23 ax.scatter(repair['longitude'], repair['latitude'], c='orange', alpha=0.5)
24 ax.scatter(non_functional['longitude'], non_functional['latitude'], c='red', alpha=0.5)
25
26 # Limiting the display area
27 ax.set_ylim(-12, 0)
28 ax.set_xlim(28, 41)
29
30 # Adding labels and title
31 ax.set_xlabel('Longitude')
32 ax.set_ylabel('Latitude')
33 ax.set_title('Well Status in Tanzania')
34
35 # Adding legend
36 ax.legend(markerscale=3, loc='upper right')
37
38 plt.show()
39
```



To further identify geographic feature relationships, we'll plot only the non-functional wells and adjust the marker size by 'gps_height'.

```
In [25]: 1 # Filter for 'non functional' wells
2 non_functional = gdf[gdf['status_group'] == 'non functional']
3
4 # Filter for Tanzania and plot
5 fig, ax = plt.subplots(figsize=(15, 15))
6 base = world[world.ADMIN == 'United Republic of Tanzania'].plot(color='white', edgecolor='black')
7
8 # Plot for 'non functional' wells with marker size based on 'gps_height'
9 # Normalize 'gps_height' for visualization
10 max_height = non_functional['gps_height'].max()
11 marker_size = (non_functional['gps_height'] / max_height) * 100
12
13 ax.scatter(non_functional['longitude'], non_functional['latitude'],
14            s=marker_size)
15 # Limiting the display area
16 ax.set_xlim(-12, 0)
17 ax.set_ylim(28, 41)
18
19 # Adding labels and title
20 ax.set_xlabel('Longitude')
21 ax.set_ylabel('Latitude')
22 ax.set_title('Non Functional Wells in Tanzania')
23
24 # Adding legend
25 ax.legend(markerscale=3, loc='upper right')
26
27 plt.show()
```



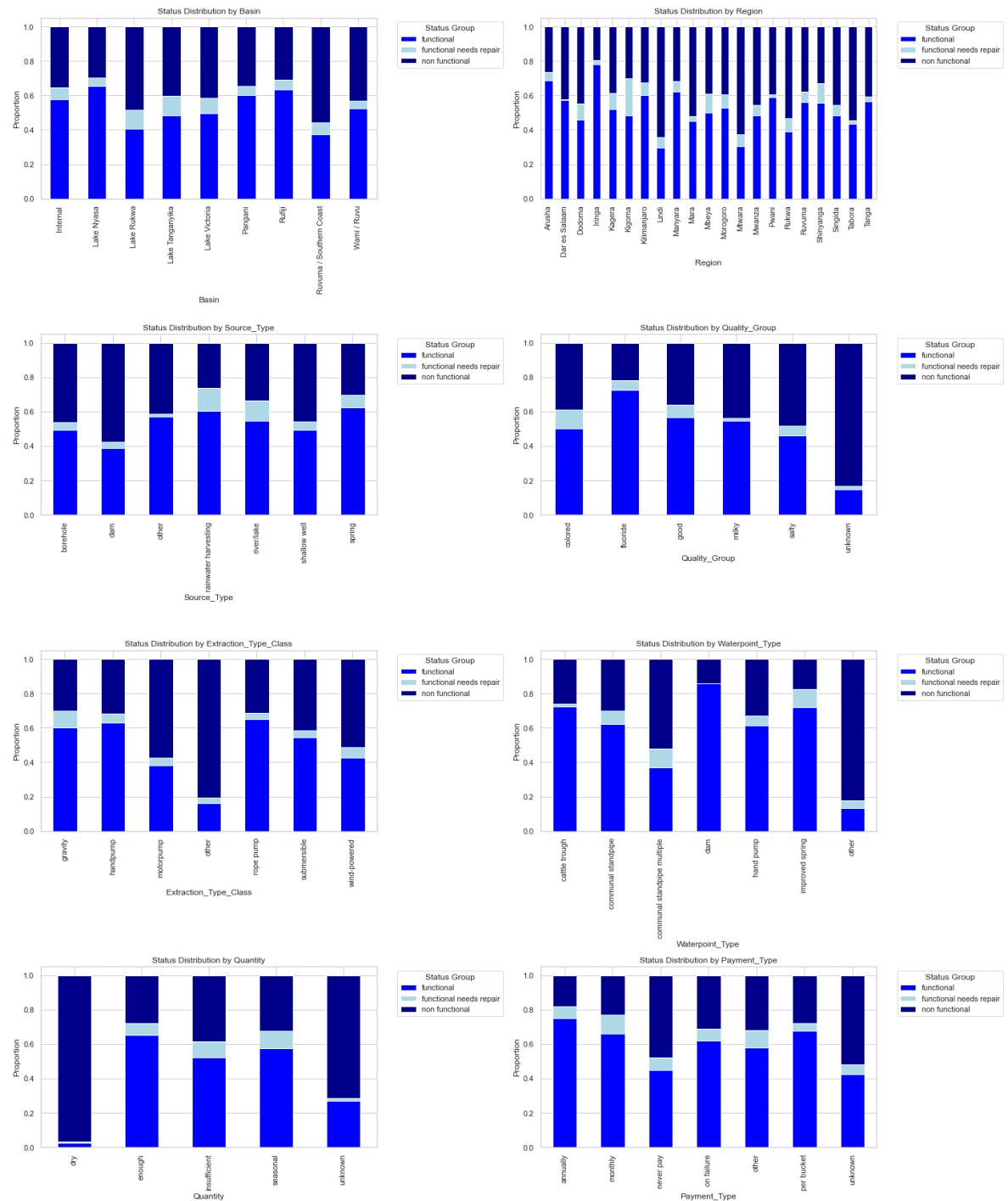
The majority of 'non functional' wells seem to be located at higher altitudes.

Categorical to Target Relationships

We'll explore the relationships between categorical variables and the target variable 'status_group'.

In [26]:

```
1 import math
2
3 categorical_vars = ['basin', 'region', 'source_type', 'quality_group']
4
5 # Define custom colors for each status
6 color_map = {'functional': 'blue', 'non functional': 'darkblue', 'no data': 'lightgray'}
7
8 # Determine the number of rows needed for the subplot (2 plots per row)
9 num_vars = len(categorical_vars)
10 num_rows = math.ceil(num_vars / 2)
11
12 # Creating subplots for each categorical variable
13 fig, axes = plt.subplots(nrows=num_rows, ncols=2, figsize=(20, num_rows * 5))
14 axes = axes.flatten() # Flatten the axes array for easy iteration
15
16 # Looping through the variables and creating a stacked bar plot for each
17 for i, var in enumerate(categorical_vars):
18     # Creating a crosstab for the variable and status_group
19     crosstab = pd.crosstab(master_data[var], master_data['status_group'])
20
21     # Creating a stacked bar plot with custom colors
22     crosstab.plot(kind='bar', stacked=True, color=[color_map[status] for status in crosstab.index])
23     axes[i].set_title(f'Status Distribution by {var.title()}')
24     axes[i].set_xlabel(var.title())
25     axes[i].set_ylabel('Proportion')
26     axes[i].legend(title='Status Group', bbox_to_anchor=(1.05, 1), borderaxespad=0)
27
28 # Adjust the layout
29 plt.tight_layout()
30 plt.show()
31
32
```



Observations

Region: Similar to basins, each region has a unique distribution of well statuses. This seems to be an indicator for well status.

Payment Type: Whether a well is paid seems to be a crucial factor. Wells that are not paid have a high number of non-functional wells.

Waterpoint Type: The method used for the population to access the water from the wells is another crucial factor. Similarly to extraction methods, waterpoint types might be more robust and less prone to failure, while others could be more complex and require frequent repairs.

Machine Learning

Baseline Model #1: Binary Target Column

Class 0 = non-functional/needs repair

Class 1 = functional

In the context of non-functional wells, focusing on recall (false negative) may be more important to ensure that most of the non-functional wells are correctly identified. We will test two separate baseline models, each with a different binary target column to see which performs best on recall.

```
In [27]: 1 # Initiating train_test_split
2 X = master_data[['waterpoint_type']]
3 y = master_data['status_binary']
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
5
6 # Encoding categorical variable
7 from sklearn.preprocessing import OneHotEncoder
8 ohe = OneHotEncoder()
9 ohe.fit(X_train)
10 X_train_encoded = ohe.transform(X_train)
11 X_test_encoded = ohe.transform(X_test)
12
13 # Plotting Log Reg transform
14 logreg = LogisticRegression(random_state=42)
15 logreg.fit(X_train_encoded, y_train)
16
17 # Checking if the target is balanced
18 y_test.value_counts(normalize=True)
19
20 y_pred = logreg.predict(X_test_encoded)
21
22 print("Classification Report:\n", classification_report(y_test, y_pred))
23 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
24
```

Classification Report:

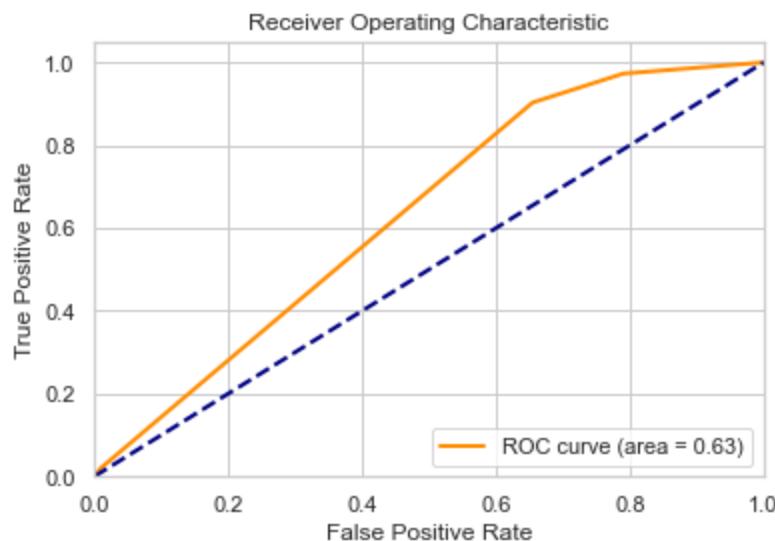
	precision	recall	f1-score	support
0	0.75	0.35	0.47	8894
1	0.62	0.90	0.73	10479
accuracy			0.65	19373
macro avg	0.69	0.62	0.60	19373
weighted avg	0.68	0.65	0.61	19373

Confusion Matrix:

```
[[3070 5824]
[1015 9464]]
```

Evaluating with ROC Curve

```
In [28]: 1 # Predict probas for the positive class
2 y_pred_proba = logreg.predict_proba(X_test_encoded)[:, 1]
3
4 # Compute AUC-ROC
5 roc_auc = roc_auc_score(y_test, y_pred_proba)
6
7 # Compute ROC curve
8 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
9
10 # Plotting ROC Curve
11 plt.figure()
12 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = 0.63)')
13 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
14 plt.xlim([0.0, 1.0])
15 plt.ylim([0.0, 1.05])
16 plt.xlabel('False Positive Rate')
17 plt.ylabel('True Positive Rate')
18 plt.title('Receiver Operating Characteristic')
19 plt.legend(loc="lower right")
20 plt.show()
```



Baseline Model #2: Reverse Binary Target Column

Class 0 = non-functional

Class 1 = functional/needs repair

```
In [29]: 1 # Initiating train_test_split
2 X = master_data[['waterpoint_type']]
3 y = master_data['status_binary_reversed']
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
5
6 # Encoding categorical variable
7 from sklearn.preprocessing import OneHotEncoder
8 ohe = OneHotEncoder()
9 ohe.fit(X_train)
10 X_train_encoded = ohe.transform(X_train)
11 X_test_encoded = ohe.transform(X_test)
12
13 # Plotting Log Reg transform
14 logreg = LogisticRegression(random_state=42)
15 logreg.fit(X_train_encoded, y_train)
16
17 # Checking if the target is balanced
18 y_test.value_counts(normalize=True)
19
20 # Predicting and evaluating the model
21 y_pred = logreg.predict(X_test_encoded)
22 print("Classification Report:\n", classification_report(y_test, y_pred))
23 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

Classification Report:

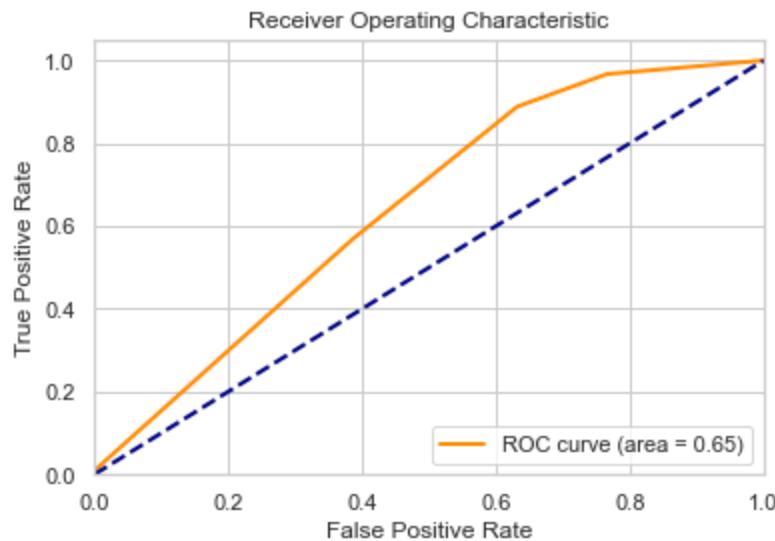
	precision	recall	f1-score	support
0	0.68	0.37	0.48	7490
1	0.69	0.89	0.78	11883
accuracy			0.69	19373
macro avg	0.68	0.63	0.63	19373
weighted avg	0.68	0.69	0.66	19373

Confusion Matrix:

```
[[ 2758  4732]
 [ 1327 10556]]
```

Evaluating with ROC Curve

```
In [30]: 1 # Predict probas for the positive class
2 y_pred_proba = logreg.predict_proba(X_test_encoded)[:, 1]
3
4 # Compute AUC-ROC
5 roc_auc = roc_auc_score(y_test, y_pred_proba)
6
7 # Compute ROC curve
8 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
9
10 # Plotting ROC Curve
11 plt.figure()
12 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = 0.65)')
13 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
14 plt.xlim([0.0, 1.0])
15 plt.ylim([0.0, 1.05])
16 plt.xlabel('False Positive Rate')
17 plt.ylabel('True Positive Rate')
18 plt.title('Receiver Operating Characteristic')
19 plt.legend(loc="lower right")
20 plt.show()
```



Baseline Comparison

Baseline 1

Precision:

- Class 0: 0.75 (Higher precision for class 0)
- Class 1: 0.62

Recall:

- Class 0: 0.35 (Lower recall for class 0)
- Class 1: 0.90 (Higher recall for class 1)

F1-Score:

- Class 0: 0.47 (Lower F1-score for class 0)
- Class 1: 0.73

▲ ----- 0.53%

Baseline 2

Precision:

- Class 0: 0.68 (Lower precision for class 0)
- Class 1: 0.69

Recall:

- Class 0: 0.37 (Slightly higher recall for class 0)
- Class 1: 0.89 (Slightly lower recall for class 1)

F1-Score:

- Class 0: 0.48 (Slightly higher F1-score for class 0)
- Class 1: 0.78

Accuracy: 69% (Higher)

Analysis

Baseline 2 shows improved overall performance, with better accuracy and a better balance in precision and recall for both classes. However, it is more prone to falsely identifying class 0 (non-functional) instances as class 1 (functional).

Baseline 1 while having higher precision for class 0 (non-functional), falls short in accurately identifying class 0 (non-functional) instances (lower recall).

Since we are more concerned with better recall, we will continue our modeling with **Baseline 2**.

Random Forest Classifier

This model can provide insights into the importance of various features in predicting well functionality. It's less likely to overfit than individual decision trees and doesn't require feature scaling.

```
In [*]: 1 X = master_data.drop(['status_binary', 'status_binary_reversed', 's
2                               'installer', 'permit', 'date_recorded', 'cons
3 y = master_data['status_binary_reversed']
4
5 # Initializing train/test split
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
7
8 # Identifying categoricals
9 categorical_cols = X_train.select_dtypes(include=['object', 'catego
10
11 # Creating a column transformer with OneHotEncoder for categoricals
12 column_transformer = ColumnTransformer(
13     transformers=[
14         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_
15     ],
16     remainder='passthrough'
17 )
18
19 # Applying the column transformer
20 X_train_encoded = column_transformer.fit_transform(X_train)
21 X_test_encoded = column_transformer.transform(X_test)
22
23 # Creating and training the Random Forest model
24 rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
25 rf_model.fit(X_train_encoded, y_train)
26
27 # Extracting feature names manually for older versions of scikit-learn
28 onehot_features = column_transformer.named_transformers_['cat'].get_
29 other_features = [col for col in X_train.columns if col not in cate
30 feature_names = np.concatenate([onehot_features, other_features])
31
32 # Predicting and evaluating the model
33 y_pred = rf_model.predict(X_test_encoded)
34 print("Classification Report:\n", classification_report(y_test, y_p
35 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
36
```

Evaluating with ROC Curve

```
In [*]: 1 # Predict probas for the positive class
2 y_pred_proba = rf_model.predict_proba(X_test_encoded)[:, 1]
3
4 # Compute AUC-ROC
5 roc_auc = roc_auc_score(y_test, y_pred_proba)
6
7 # Compute ROC curve
8 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
9
10 # Plotting ROC Curve
11 plt.figure()
12 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area
13 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
14 plt.xlim([0.0, 1.0])
15 plt.ylim([0.0, 1.05])
16 plt.xlabel('False Positive Rate')
17 plt.ylabel('True Positive Rate')
18 plt.title('Receiver Operating Characteristic')
19 plt.legend(loc="lower right")
20 plt.show()
```

Evaluating Feature Importance

```
In [*]: 1 # Extracting and displaying feature importances
2 importances = rf_model.feature_importances_
3 importance_df = pd.DataFrame({'Feature': feature_names, 'Importance':
4 top_features = importance_df.sort_values(by='Importance', ascending=
5 print(top_features)
```

Tuning Random Forest Classifier

- SMOTE for oversampling the minority class or adjusting class weights in the model.
- Hyperparameter tuning of the Random Forest model
- RandomizedSearchCV will randomly sample 10 combos of parameters and use 3-fold cross-validation. This will reduce run time compared to GridSearchCV

```
In [*]: 1 from imblearn.over_sampling import SMOTE
2 from sklearn.model_selection import RandomizedSearchCV
3
4 # Handling class imbalance with SMOTE
5 smote = SMOTE()
6 X_train_resampled, y_train_resampled = smote.fit_resample(X_train_encoded, y_train)
7
8 # Define the hyperparameter grid
9 param_grid = {
10     'n_estimators': [100, 200, 300],
11     'max_depth': [10, 20, 30],
12     'min_samples_split': [2, 5, 10],
13     'min_samples_leaf': [1, 2, 4],
14     'max_features': ['sqrt', 'log2']
15 }
16
17 # Hyperparameter tuning with Randomized Search
18 random_search = RandomizedSearchCV(
19     RandomForestClassifier(random_state=42),
20     param_grid,
21     n_iter=10,
22     cv=3,
23     scoring='recall',
24     n_jobs=-1
25 )
26
27 random_search.fit(X_train_resampled, y_train_resampled)
28
29 # Get the best model
30 best_model = random_search.best_estimator_
31
32 # Re-train and evaluate the model with the best params
33 best_model.fit(X_train_resampled, y_train_resampled)
34 y_pred = best_model.predict(X_test_encoded)
35
36 # Evaluate the model
37 print("Classification Report:\n", classification_report(y_test, y_pred))
38 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
39
```

Evaluating with ROC Curve

```
In [*]: 1 # Predict probas for the positive class
2 y_pred_proba = rf_model.predict_proba(X_test_encoded)[:, 1]
3
4 # Compute AUC-ROC
5 roc_auc = roc_auc_score(y_test, y_pred_proba)
6
7 # Compute ROC curve
8 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
9
10 # Plotting ROC Curve
11 plt.figure()
12 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area
13 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
14 plt.xlim([0.0, 1.0])
15 plt.ylim([0.0, 1.05])
16 plt.xlabel('False Positive Rate')
17 plt.ylabel('True Positive Rate')
18 plt.title('Receiver Operating Characteristic')
19 plt.legend(loc="lower right")
20 plt.show()
```

Comparison Interpretations of Random Forest Models

Class 0 = non-functional/needs repair

Class 1 = functional

Precision:

- **Original Model:** Precision is 0.83 for both classes.
- **Tuned Model:** Precision increased to 0.85 for class 0 but decreased to 0.79 for class 1.

Recall:

- **Original Model:** Recall is 0.77 for class 0 and 0.90 for class 1.
- **Tuned Model:** Recall decreased to 0.62 for class 0 but increased to 0.93 for class 1.

F1-Score:

- **Original Model:** F1-scores are 0.80 (class 0) and 0.88 (class 1).
- **Tuned Model:** F1-scores are 0.72 (class 0) and 0.86 (class 1).

Accuracy:

- **Original Model:** Overall accuracy is 0.85.
- **Tuned Model:** Overall accuracy decreased to 0.81.

Macro and Weighted Averages:

- **Original Model:** Both macro and weighted averages are around 0.85.

- **Tuned Model:** Both macro and weighted averages are around 0.77 - 0.81.

AUC-ROC Score:

- **Original Model:** The area under the curve is 0.91, which is high. This means the model can effectively distinguish between the positive class (class 1) and the negative class (class 0).
- **Tuned Model:** The area under the curve is also 0.91.

Analysis of Comparison:

- Tuning the model appears to have made it more biased towards class 1, improving its ability to detect class 1 instances but worsening its performance for class 0 (higher false positives).
- The original model is more balanced in terms of precision and recall across both classes.
- The tuned model has a lower overall accuracy compared to the original model.

Logistic Regression Model: Most Important Features

```
In [*]: 1 # Selecting most important features
2 features = ['region', 'quantity', 'gps_height', 'waterpoint_type',
3 X = master_data[features]
4 y = master_data['status_binary_reversed']
5
6 # Initializing train/test split
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
8
9 # Identifying categorical and numerical
10 categorical_cols = ['quantity', 'region', 'waterpoint_type', 'extra_freshwater_source']
11 numerical_cols = ['gps_height']
12
13 # Creating a column transformer with OneHotEncoder and StandardScaler
14 preprocessor = ColumnTransformer(
15     transformers=[
16         ('num', StandardScaler(), numerical_cols),
17         ('cat', OneHotEncoder(), categorical_cols)
18     ]
19 )
20
21 # Creating a pipeline with preprocessing and log reg model
22 model_pipeline = Pipeline([
23     ('preprocessor', preprocessor),
24     ('classifier', LogisticRegression(solver='saga', max_iter=1000))
25 ])
26
27 # Training the model
28 model_pipeline.fit(X_train, y_train)
29
30 # Predicting and evaluating the model
31 y_pred = model_pipeline.predict(X_test)
32 print("Classification Report:\n", classification_report(y_test, y_pred))
33 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
34
```

```
In [*]: 1 # Calculate probas, ROC curve, and AUC for log reg
2 logreg_probs = model_pipeline.predict_proba(X_test)[:, 1]
3 fpr_logreg, tpr_logreg, _ = roc_curve(y_test, logreg_probs)
4 roc_auc_logreg = auc(fpr_logreg, tpr_logreg)
5
6
7 # Plotting ROC Curve
8 plt.figure()
9 plt.plot(fpr_logreg, tpr_logreg, color='darkorange', lw=2, label='Logistic Regression')
10 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
11 plt.xlim([0.0, 1.0])
12 plt.ylim([0.0, 1.05])
13 plt.xlabel('False Positive Rate')
14 plt.ylabel('True Positive Rate')
15 plt.title('Receiver Operating Characteristic')
16 plt.legend(loc="lower right")
17 plt.show()
```

Interpretation

Precision:

- For class 0: 83% precision means that out of all instances predicted as class 0, 83% were actually class 0.
- For class 1: 76% precision indicates that out of all instances predicted as class 1, 76% were actually class 1.

Recall:

- For class 0: The recall of 54% is moderate, meaning the model correctly identifies 54% of the actual class 0 instances.
- For class 1: A high recall of 93% indicates the model is very effective at identifying class 1 instances.

F1-Score:

- For class 0: The F1-score of 0.65 suggests a balance between precision and recall for class 0, but more weighted towards precision.
- For class 1: The F1-score of 0.84 indicates a strong balance between precision and recall for class 1, favoring recall.

Accuracy:

- The overall accuracy of 78% indicates that the model correctly predicts the class for 78% of all instances.

Macro and Weighted Averages:

- Macro average treats both classes equally, showing an average precision of 79%, recall of 73%, and F1-score of 74%.
- Weighted average considers class imbalance, showing slightly higher precision and recall, indicating better performance on the more prevalent class 1.

ROC Score:

- An ROC score of 0.82 suggests a good ability of the model to distinguish between the two classes. It indicates a favorable balance between the true positive rate and false positive rate across different thresholds.

Insights:

- The model performs well overall, especially in predicting class 1, which is indicated by the high recall and F1-score for class 1.
- Model is less effective in correctly identifying class 0 instances, as evidenced by the lower recall for class 0.
- The relatively high number of false positives for class 0 (3162) indicates that the model often misclassifies class 1 instances as class 0.

Key Findings

- **Geographic Indicators:** Including region and altitude, geographic features are 21% MORE influential in identifying non-functional wells than other features.
- **Region:** Mbeya, Morogoro, and Kilimanjaro have the highest rates of non-functional wells. Altitude may play an important role in water source access.
- **Type of Wells:** Communal Standpipe wells are most likely to be a functional well. Other well types have the highest percentage of non-functional wells at 81.38%. Multi Communal Handpipe wells have the second highest percentage of non-functional wells at 53.85%.
- **Payment Type:** Whether a well is paid seems to be a crucial factor. Wells that are not paid have a high number of non-functional wells.
- **Random Forest Classifier:** Our best performing model gave us actionable insights into feature importance and effectively minimized false-negatives.

Conclusion

In this project, we have unearthed critical insights to steer the strategic decision-making of the Tanzanian Government. This includes pinpointing the precise types of wells that warrant prioritized construction efforts, as well as identifying the specific regions that should receive initial focus and substantial investment in well infrastructure.

Next Steps

- **Investigate Additional Features:** Concentrating on geographical indicators like climate, population, and amount of water available in the area.
- **Time-Series Analysis:** Further consideration of the well ages should be analyzed to predict the average lifetime of more robust well structures.
- **Repairs:** Local governments should look at what type of water wells are needing repairs, and the severity of those repairs, to fine-tune non-functional indicators.

Sources

- [Driven Data - Tanzanian Water Wells](https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/page/23/) (<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/page/23/>)
 - [Labels](https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/) (<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/>)
 - [Values](https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/) (<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/>)
- [The World Bank](https://www.worldbank.org/en/news/feature/2022/06/14/the-road-to-better-services-is-paved-with-strong-delivery-institutions-a-rural-water-story) (<https://www.worldbank.org/en/news/feature/2022/06/14/the-road-to-better-services-is-paved-with-strong-delivery-institutions-a-rural-water-story>)
- [Groundwater Wells](https://www.usgs.gov/special-topics/water-science-school/science/groundwater-wells) (<https://www.usgs.gov/special-topics/water-science-school/science/groundwater-wells>)
- [Detection of Non-Function Bore Wells Using Maching Learning Algorithms](https://www.researchgate.net/publication/349446319_Detection_of_Non-functional_Bore_wells_Using_Machine_Learning_Algorithms) ([https://www.researchgate.net/publication/349446319 Detection of Non-functional Bore wells Using Machine Learning Algorithms](https://www.researchgate.net/publication/349446319_Detection_of_Non-functional_Bore_wells_Using_Machine_Learning_Algorithms))

About Us

Goknur Kaya

- Github Lead
- Email: goknurkaya@gmail.com (<mailto:goknurkaya@gmail.com>)
- github.com/goknurk

Kari Primiano

- Tech Lead
- Email: kkprim@gmail.com (<mailto:kkprim@gmail.com>)
- github.com/kkprim

Mytreyi Abburu

- Presentation Lead
- Email: abburumk@gmail.com (<mailto:abburumk@gmail.com>)
- github.com/myt-hue