

Tanzanian Water Wells

Problem Description



Water Aid is an NGO based in the United Kingdom that works on access to clean water around the world. They consider access to clean water, decent toilets and good hygiene as basic human rights. For over 30 years, they have been working in partnership to improve access to these three essentials through a combination of programmatic and policy work.

Water Aid works in several countries around the globe, including Tanzania. According to the World Sector Report (2019) around 60% of Tanzanians have access to improved water, but the degree of water access, and the water quality and quantity, varies. Drought, landscape change, and the amplifying effects of climate change are straining existing surface water supplies.

Water Aid is launching a program to repair non-functioning wells in the cross country shared water basins of Eastern Africa. The status of the wells is not clearly recorded in countries surrounding Tanzania. Identifying non-functioning wells, securing funding, and traveling to these rural locations to repair wells is both time and resource intensive. They need a predictive model that accurately identifies which wells are not functioning to reduce cost and ensure they are using their resources wisely. They also need to identify a specific water basin to begin their work.

Goals

There are three data science goals to address Water Aid's need for accurately identifying non-functioning wells:

1. Using an iterative process, build a predictive machine learning model based on existing water well data to accurately classify non-functioning wells.
2. Deliver two recommendations to Water Aid: a specific transboundary water basin to begin their operations, and one feature characteristic of the wells in this basin that will

lead to higher chance of identifying non-functioning wells.

Load Packages and Data

```
In [368]: 1 import pandas as pd
          2 import numpy as np
          3 from matplotlib import pyplot as plt
          4 import seaborn as sns
          5
          6 %matplotlib inline
          7
          8 from sklearn.model_selection import train_test_split, GridSearchCV
          9 from sklearn.pipeline import Pipeline
         10 from sklearn.preprocessing import StandardScaler, OneHotEncoder, F
         11 from sklearn.impute import SimpleImputer
         12 from sklearn.compose import ColumnTransformer
         13 from sklearn.linear_model import LogisticRegression
         14 from sklearn.tree import DecisionTreeClassifier
         15 from sklearn.ensemble import RandomForestClassifier, GradientBoost
         16 from sklearn.metrics import plot_confusion_matrix, recall_score,\
         17     accuracy_score, precision_score, f1_score
         18
         19 from imblearn.over_sampling import SMOTE
         20 from imblearn.pipeline import Pipeline as ImPipeline
```

```
In [369]: 1 # Load the predictor data
          2
          3 wells = pd.read_csv('training_set_values.csv')
          4 wells.head()
```

Out [369]:

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latitude	w
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	
4	19728	0.0	2011-07-13	Action In A	0	Artisan	31.130847	-1.825359	

5 rows × 10 columns

```
In [370]: 1 # Load the target data
          2
          3 target = pd.read_csv('training_set_labels.csv')
          4
          5 target.head()
```

Out [370]:

	id	status_group
0	69572	functional
1	8776	functional
2	34310	functional
3	67743	non functional
4	19728	functional

1. Exploratory Data Analysis

Get a sense of the big picture for the dataset. Prepare the data for further analysis. Gain an understanding of the variables, or predictors in this case. Study the relationship between variables. Make plan for initial model.

```
In [371]: 1 # Identify size of dataset
          2
          3 print("Records for wells:", wells.shape)
          4 print()
          5 print("Records for target:", target.shape)
```

Records for wells: (59400, 40)

Records for target: (59400, 2)

```
In [372]: 1 # Identify datatypes and record amount for each predictor
          2
          3 wells.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                               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
dtypes: float64(3), int64(7), object(30)
memory usage: 18.1+ MB
```

```
In [373]: 1 # Examine numerical predictors mean, min, max
          2
          3 wells.describe()
```

Out[373]:

	id	amount_tsh	gps_height	longitude	latitude	num_private
count	59400.000000	59400.000000	59400.000000	59400.000000	5.940000e+04	59400.000000
mean	37115.131768	317.650385	668.297239	34.077427	-5.706033e+00	0.474141
std	21453.128371	2997.574558	693.116350	6.567432	2.946019e+00	12.236230
min	0.000000	0.000000	-90.000000	0.000000	-1.164944e+01	0.000000
25%	18519.750000	0.000000	0.000000	33.090347	-8.540621e+00	0.000000
50%	37061.500000	0.000000	369.000000	34.908743	-5.021597e+00	0.000000
75%	55656.500000	20.000000	1319.250000	37.178387	-3.326156e+00	0.000000
max	74247.000000	350000.000000	2770.000000	40.345193	-2.000000e-08	1776.000000

```
In [374]: 1 # Missing data total
          2 wells.isna().sum().sum()
```

Out[374]: 46094

```
In [375]: 1 # Missing data by predictor
          2
          3 wells.isna().sum()
```

```
Out[375]: id                                0
          amount_tsh                        0
          date_recorded                     0
          funder                            3635
          gps_height                        0
          installer                         3655
          longitude                        0
          latitude                         0
          wpt_name                         0
          num_private                      0
          basin                            0
          subvillage                       371
          region                          0
          region_code                     0
          district_code                   0
          lga                             0
          ward                            0
          population                      0
          public_meeting                  3334
          recorded_by                     0
          scheme_management               3877
          scheme_name                    28166
          permit                         3056
          construction_year               0
          extraction_type                 0
          extraction_type_group           0
          extraction_type_class           0
          management                     0
          management_group                0
          payment                        0
          payment_type                   0
          water_quality                   0
          quality_group                   0
          quantity                       0
          quantity_group                  0
          source                         0
          source_type                     0
          source_class                    0
          waterpoint_type                 0
          waterpoint_type_group           0
          dtype: int64
```

```
In [376]: 1 # Data missing for target
          2 target.isna().sum()
```

```
Out[376]: id                                0
          status_group                      0
          dtype: int64
```

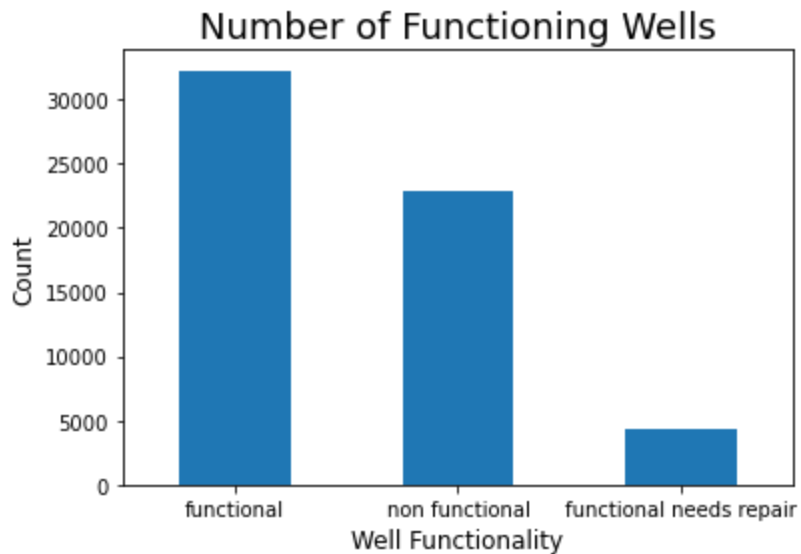
```
In [377]: 1 # Examine value counts for the target, consider imbalance in target
          2 target['status_group'].value_counts()
```

```
Out[377]: functional                32259
          non functional            22824
          functional needs repair   4317
          Name: status_group, dtype: int64
```

```
In [378]: 1 # Percentage makeup of target values
          2 print("Functional percentage:", round(32259/59400*100, 2))
          3 print("Non functional percentage:", round(22824/59400*100, 2))
          4 print("Functional needs repair percentage:", round(4317/59400*100,
```

```
Functional percentage: 54.31
Non functional percentage: 38.42
Functional needs repair percentage: 7.27
```

```
In [379]: 1 # Visually plot target variable counts
          2
          3 target.status_group.value_counts().plot(kind="bar")
          4 plt.title("Number of Functioning Wells", fontsize= 18)
          5 plt.xlabel("Well Functionality", fontsize = 12)
          6 plt.xticks(rotation=0)
          7 plt.ylabel("Count", fontsize = 12)
          8 plt.show();
          9
         10 plt.savefig('Number of Functioning Wells')
```



<Figure size 432x288 with 0 Axes>

```
In [380]: 1 # Identify unique values per column  
          2 print(wells.nunique())
```

```
id                59400  
amount_tsh        98  
date_recorded     356  
funder            1897  
gps_height        2428  
installer         2145  
longitude         57516  
latitude          57517  
wpt_name          37400  
num_private       65  
basin             9  
subvillage        19287  
region            21  
region_code       27  
district_code     20  
lga               125  
ward              2092  
population        1049  
public_meeting    2  
recorded_by       1  
scheme_management 12  
scheme_name       2696  
permit            2  
construction_year 55  
extraction_type    18  
extraction_type_group 13  
extraction_type_class 7  
management         12  
management_group   5  
payment            7  
payment_type       7  
water_quality      8  
quality_group      6  
quantity           5  
quantity_group     5  
source             10  
source_type        7  
source_class       3  
waterpoint_type    7  
waterpoint_type_group 6  
dtype: int64
```



```
In [381]: 1 # Concatenate preds and target for heatmap
          2
          3 df = pd.concat([wells, target], axis =1)
          4
          5 df = df.loc[:,~df.columns.duplicated()].copy()
          6
          7 df.head()
```

Out [381]:

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latitude	v
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	
4	19728	0.0	2011-07-13	Action In A	0	Artisan	31.130847	-1.825359	

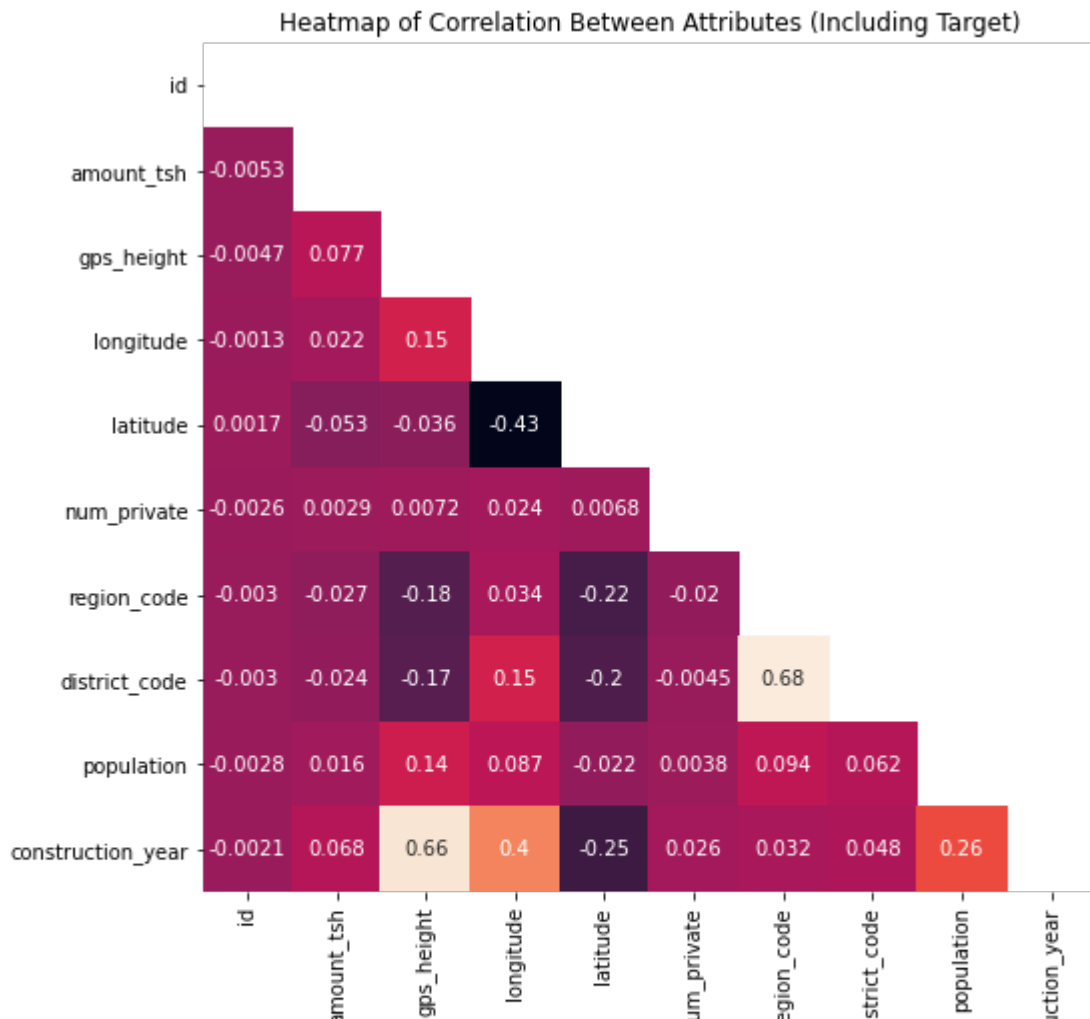
5 rows × 41 columns

Correlation of numeric data

```

In [382]: 1 # Create a heatmap to examine the correlational coefficients
          2
          3
          4 corr = df.corr()
          5
          6 # Set up figure and axes
          7 fig, ax = plt.subplots(figsize=(8, 12))
          8
          9 # Plot a heatmap of the correlations
         10
         11 sns.heatmap(
         12     data=corr,
         13     mask=np.triu(np.ones_like(corr, dtype=bool)),
         14     ax=ax,
         15     annot=True,
         16     # Customizes colorbar appearance
         17     cbar_kws={"label": "Correlation", "orientation": "horizontal",
         18             }
         19 )
         20
         21 # Customize the plot appearance
         22
         23
         24 ax.set_title("Heatmap of Correlation Between Attributes (Including
         25

```

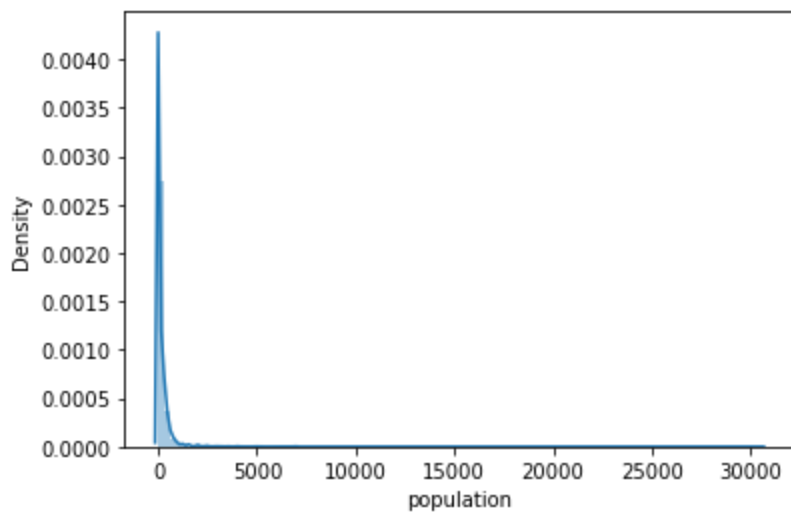


```

In [383]: 1 # Check distribution for numeric data
          2 import warnings
          3 warnings.filterwarnings('ignore')
          4
          5 print(sns.distplot(wells.population, bins = 100))
          6
          7

```

AxesSubplot(0.125,0.125;0.775x0.755)

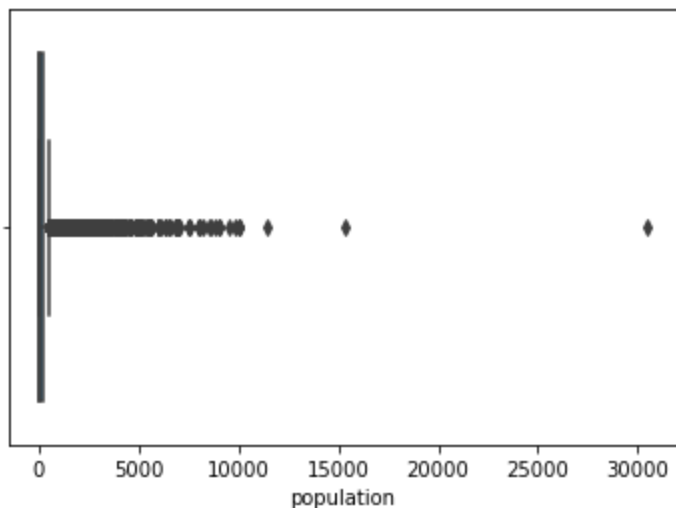


```

In [384]: 1 # Identify range of population and any outliers
          2
          3 sns.boxplot(wells.population)
          4

```

Out[384]: <AxesSubplot:xlabel='population'>



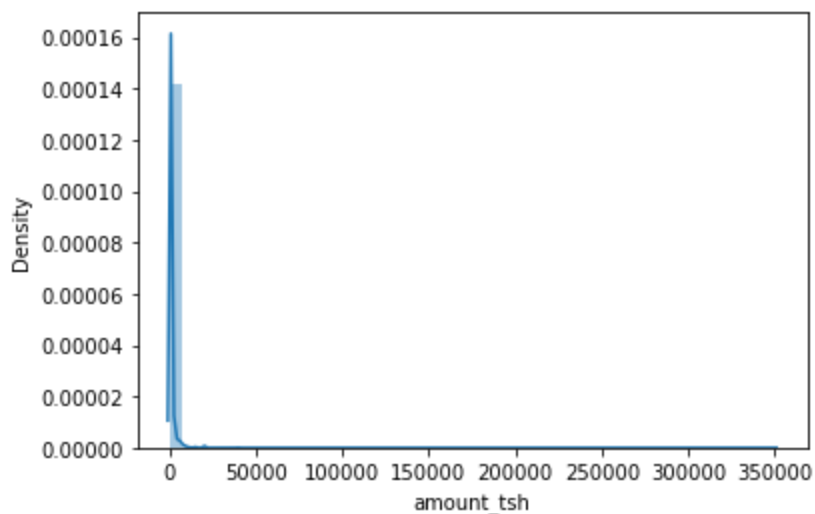
```
In [385]: 1 # Identify population counts
          2 print(wells.population.value_counts())
          3 print(wells.population.nunique())
```

```
0      21381
1       7025
200     1940
150     1892
250     1681
...
3241         1
1960         1
1685         1
2248         1
1439         1
Name: population, Length: 1049, dtype: int64
1049
```

Total Static Head data

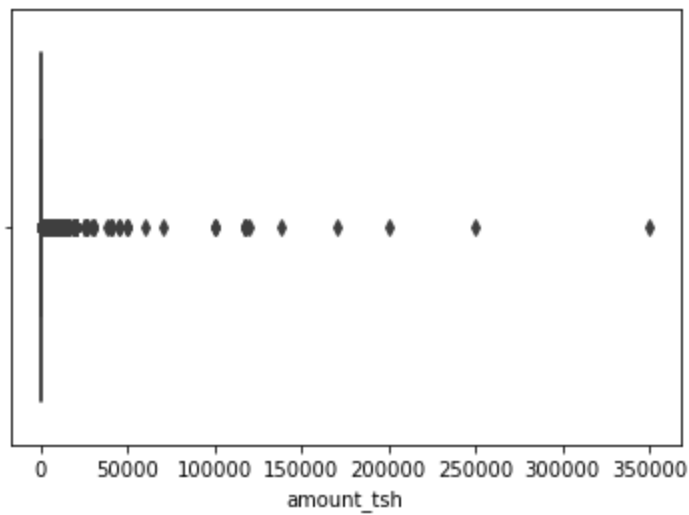
```
In [386]: 1 # Plot and describe total static head
          2
          3 sns.distplot(wells.amount_tsh)
          4
          5 print(wells.amount_tsh.describe())
          6
```

```
count      59400.000000
mean        317.650385
std         2997.574558
min           0.000000
25%           0.000000
50%           0.000000
75%          20.000000
max      350000.000000
Name: amount_tsh, dtype: float64
```



```
In [387]: 1 # Identify range and outliers of total static head
          2
          3 sns.boxplot(wells.amount_tsh)
```

```
Out[387]: <AxesSubplot:xlabel='amount_tsh'>
```

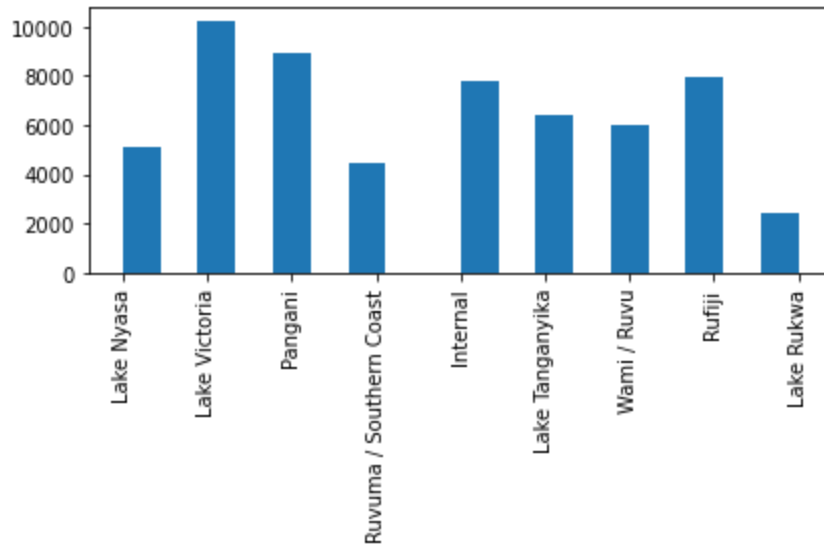


```
In [388]: 1 # Identify how many wells do not have static head
          2
          3 print(wells.amount_tsh.value_counts())
```

```
0.0          41639
500.0         3102
50.0          2472
1000.0        1488
20.0          1463
...
8500.0         1
6300.0         1
220.0          1
138000.0       1
12.0           1
Name: amount_tsh, Length: 98, dtype: int64
```

Geographic Data

```
In [389]: 1 # plot basins
2 fig, axs = plt.subplots(1, 1,
3                       figsize=(6, 4),
4                       tight_layout = True)
5
6 axs.hist(wells.basin, bins = 18)
7 plt.xticks(rotation = 90)
8
9 plt.show()
```



```
In [390]: 1 # How many wells are in each basin
2
3 wells.basin.value_counts()
```

```
Out[390]: Lake Victoria      10248
Pangani          8940
Rufiji           7976
Internal         7785
Lake Tanganyika  6432
Wami / Ruvu     5987
Lake Nyasa       5085
Ruvuma / Southern Coast  4493
Lake Rukwa       2454
Name: basin, dtype: int64
```

```
In [391]: 1  
          2 wells.region.value_counts()
```

```
Out[391]: Iringa          5294  
          Shinyanga      4982  
          Mbeya          4639  
          Kilimanjaro    4379  
          Morogoro       4006  
          Arusha         3350  
          Kagera         3316  
          Mwanza         3102  
          Kigoma         2816  
          Ruvuma         2640  
          Pwani          2635  
          Tanga          2547  
          Dodoma         2201  
          Singida        2093  
          Mara           1969  
          Tabora         1959  
          Rukwa          1808  
          Mtwara         1730  
          Manyara       1583  
          Lindi          1546  
          Dar es Salaam   805  
          Name: region, dtype: int64
```

```
In [392]: 1 wells.district_code.value_counts()
```

```
Out[392]: 1      12203  
          2      11173  
          3       9998  
          4       8999  
          5       4356  
          6       4074  
          7       3343  
          8       1043  
          30       995  
          33       874  
          53       745  
          43       505  
          13       391  
          23       293  
          63       195  
          62       109  
          60        63  
          0         23  
          80         12  
          67          6  
          Name: district_code, dtype: int64
```

```
In [393]: 1 wells.region_code.value_counts()
```

```
Out[393]: 11    5300
          17    5011
          12    4639
           3    4379
           5    4040
          18    3324
          19    3047
           2    3024
          16    2816
          10    2640
           4    2513
           1    2201
          13    2093
          14    1979
          20    1969
          15    1808
           6    1609
          21    1583
          80    1238
          60    1025
          90     917
           7     805
          99     423
           9     390
          24     326
           8     300
          40        1
          Name: region_code, dtype: int64
```

Water attributes

```
In [394]: 1 wells.water_quality.value_counts()
```

```
Out[394]: soft                50818
          salty                4856
          unknown             1876
          milky                804
          coloured            490
          salty abandoned      339
          fluoride             200
          fluoride abandoned    17
          Name: water_quality, dtype: int64
```



```
In [395]: 1 wells.quality_group.value_counts()
```

```
Out[395]: good          50818
          salty         5195
          unknown      1876
          milky         804
          colored       490
          fluoride      217
          Name: quality_group, dtype: int64
```

```
In [396]: 1 wells.quantity.value_counts()
```

```
Out[396]: enough        33186
          insufficient   15129
          dry            6246
          seasonal      4050
          unknown        789
          Name: quantity, dtype: int64
```

```
In [397]: 1 wells.quantity_group.value_counts()
```

```
Out[397]: enough        33186
          insufficient   15129
          dry            6246
          seasonal      4050
          unknown        789
          Name: quantity_group, dtype: int64
```

```
In [398]: 1 wells.scheme_name.value_counts()
```

```
Out[398]: K              682
          None           644
          Borehole       546
          Chalinze wate  405
          M              400
          ...
          BFFS water supplying  1
          Lwamgasa           1
          BL Embokoi         1
          Ung'oro Water supply  1
          Ugalla water supply  1
          Name: scheme_name, Length: 2696, dtype: int64
```

```
In [399]: 1 wells.scheme_management.value_counts()
```

```
Out[399]: 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
          None           1
          Name: scheme_management, dtype: int64
```

```
In [400]: 1 wells.extraction_type.value_counts()
```

```
Out[400]: 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: extraction_type, dtype: int64
```

```
In [401]: 1 wells.extraction_type_group.value_counts()
```

```
Out[401]: 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: extraction_type_group, dtype: int64
```

```
In [402]: 1 wells.extraction_type_class.value_counts()
```

```
Out[402]: gravity          26780
handpump          16456
other             6430
submersible       6179
motorpump         2987
rope pump         451
wind-powered      117
Name: extraction_type_class, dtype: int64
```

```
In [403]: 1 wells.source.value_counts()
```

```
Out[403]: 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: source, dtype: int64
```

```
In [404]: 1 wells.source_type.value_counts()
```

```
Out[404]: spring          17021
shallow well          16824
borehole             11949
river/lake           10377
rainwater harvesting   2295
dam                   656
other                 278
Name: source_type, dtype: int64
```

```
In [405]: 1 wells.waterpoint_type_group.value_counts()
```

```
Out[405]: communal standpipe  34625
hand pump                    17488
other                        6380
improved spring              784
cattle trough                116
dam                           7
Name: waterpoint_type_group, dtype: int64
```

```
In [406]: 1 wells.source_class.value_counts()
```

```
Out[406]: groundwater    45794
surface                 13328
unknown                  278
Name: source_class, dtype: int64
```

```
In [407]: 1 wells.waterpoint_type.value_counts()
```

```
Out[407]: communal standpipe      28522  
hand pump      17488  
other      6380  
communal standpipe multiple      6103  
improved spring      784  
cattle trough      116  
dam      7  
Name: waterpoint_type, dtype: int64
```

Organizational attributes

```
In [408]: 1 wells.funder.value_counts()
```

```
Out[408]: Government Of Tanzania      9084  
Danida      3114  
Hesawa      2202  
Rwssp      1374  
World Bank      1349  
...  
Mungaya      1  
Fpct Mulala      1  
Boma Saving      1  
Care Int      1  
D Ct      1  
Name: funder, Length: 1897, dtype: int64
```

```
In [409]: 1 wells.num_private.value_counts()
```

```
Out[409]: 0      58643  
6      81  
1      73  
5      46  
8      46  
...  
180      1  
213      1  
23      1  
55      1  
94      1  
Name: num_private, Length: 65, dtype: int64
```

```
In [410]: 1 wells.permit.value_counts()
```

```
Out[410]: True      38852  
False      17492  
Name: permit, dtype: int64
```

```
In [411]: 1 wells.management.value_counts()
```

```
Out[411]: 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: management, dtype: int64
```

```
In [412]: 1 wells.management_group.value_counts()
```

```
Out[412]: user-group    52490
          commercial     3638
          parastatal     1768
          other          943
          unknown        561
          Name: management_group, dtype: int64
```

```
In [413]: 1 wells.payment.value_counts()
```

```
Out[413]: never pay      25348
          pay per bucket  8985
          pay monthly     8300
          unknown         8157
          pay when scheme fails 3914
          pay annually     3642
          other           1054
          Name: payment, dtype: int64
```

```
In [414]: 1 wells.payment_type.value_counts()
```

```
Out[414]: never pay      25348
          per bucket     8985
          monthly        8300
          unknown        8157
          on failure      3914
          annually        3642
          other           1054
          Name: payment_type, dtype: int64
```

```
In [415]: 1 with pd.option_context('display.max_rows', 5, 'display.max_columns'
2          display(wells[1000:1020]))
```

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latit
1000	47384	250.0	2013-02-14	Oxfam	1409	OXFAM	30.105401	-4.367
1001	11570	0.0	2012-10-12	Resolute Mining	0	Consulting Engineer	33.210098	-4.049
...
1018	41433	0.0	2011-03-05	Government Of Tanzania	1307	DWE	38.325050	-4.464
1019	21810	0.0	2013-01-17	Bulyahunlu Gold Mine	0	Bulyahunlu Gold Mine	32.370100	-3.287

20 rows × 40 columns

2. Preprocess data, Initial Model

Redundant data columns where the data is included in other columns that contain more expansive information should be dropped: water attributes, geographic attributes, include water include regional columns, water extraction and source types.

Drop columns that do not contribute to the model. These include water id, names of the waterpoint, names of subvillages.

Make plan for missing categorical and numeric data.

```
In [416]: 1 # Drop redundant data columns and columns that do not contribute t
2 wells.drop(columns = ['id', 'wpt_name', 'region', 'recorded_by', '
3               'scheme_management', 'extraction_type_group', 'payment_
4               'quality_group', 'quantity_group', 'source_type', 'wate
5
```

In [417]: 1 wells.columns

Out[417]: Index(['amount_tsh', 'date_recorded', 'funder', 'gps_height', 'installer',
'longitude', 'latitude', 'num_private', 'basin', 'region_code',
'district_code', 'lga', 'ward', 'population', 'public_meeting',
'permit', 'construction_year', 'extraction_type',
'extraction_type_class', 'management', 'management_group', 'payment',
'water_quality', 'quantity', 'source', 'source_class',
'waterpoint_type_group'],
dtype='object')

In [418]: 1 *# Check missing data*
2 wells.isna().sum()

Out[418]: amount_tsh 0
date_recorded 0
funder 3635
gps_height 0
installer 3655
longitude 0
latitude 0
num_private 0
basin 0
region_code 0
district_code 0
lga 0
ward 0
population 0
public_meeting 3334
permit 3056
construction_year 0
extraction_type 0
extraction_type_class 0
management 0
management_group 0
payment 0
water_quality 0
quantity 0
source 0
source_class 0
waterpoint_type_group 0
dtype: int64

```
In [419]: 1 # Replace Nan in public_meeting and permit as False
          2
          3 wells['public_meeting'] = wells['public_meeting'].fillna('False').
          4 wells.public_meeting.head()
```

```
Out[419]: 0    True
          1    True
          2    True
          3    True
          4    True
          Name: public_meeting, dtype: bool
```

```
In [420]: 1 # replace missing permit data as False
          2 wells['permit'] = wells['permit'].fillna('False').astype('bool')
          3 wells.permit.head()
```

```
Out[420]: 0    False
          1     True
          2     True
          3     True
          4     True
          Name: permit, dtype: bool
```

```
In [421]: 1 # Convert "date_recorded" to month_recorded
          2
          3 import datetime
          4
          5 wells['date_recorded'] = pd.to_datetime(wells['date_recorded'])
          6 wells['month_recorded'] = wells['date_recorded'].dt.month
          7 wells['month_recorded']
```

```
Out[421]: 0         3
          1         3
          2         2
          3         1
          4         7
          ..
          59395     5
          59396     5
          59397     4
          59398     3
          59399     3
          Name: month_recorded, Length: 59400, dtype: int64
```

```
In [422]: 1 wells.drop('date_recorded', axis = 1, inplace = True)
```

Initial Model - Logistic Regression

Use a Logistic Regression model in a pipeline for initial model results.


```
In [423]: 1 # Assign the predictors and target
          2 X = wells
          3 y = target['status_group']
```

```
In [424]: 1 # Perform a train test split
          2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [425]: 1 X_train.columns
```

```
Out[425]: Index(['amount_tsh', 'funder', 'gps_height', 'installer', 'longitude',
                'latitude', 'num_private', 'basin', 'region_code', 'district_code',
                'lga', 'ward', 'population', 'public_meeting', 'permit',
                'construction_year', 'extraction_type', 'extraction_type_class',
                'management', 'management_group', 'payment', 'water_quality',
                'quantity', 'source', 'source_class', 'waterpoint_type_group',
                'month_recorded'],
                dtype='object')
```

```
In [426]: 1 # Examine data types and record counts
          2 X_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 44550 entries, 24947 to 56422
Data columns (total 27 columns):
 #   Column                                  Non-Null Count  Dtype
---  -
 0   amount_tsh                             44550 non-null  float64
 1   funder                                 41859 non-null  object
 2   gps_height                             44550 non-null  int64
 3   installer                              41850 non-null  object
 4   longitude                              44550 non-null  float64
 5   latitude                               44550 non-null  float64
 6   num_private                            44550 non-null  int64
 7   basin                                  44550 non-null  object
 8   region_code                            44550 non-null  int64
 9   district_code                          44550 non-null  int64
10   lga                                     44550 non-null  object
11   ward                                   44550 non-null  object
12   population                             44550 non-null  int64
13   public_meeting                         44550 non-null  bool
14   permit                                 44550 non-null  bool
15   construction_year                     44550 non-null  int64
16   extraction_type                        44550 non-null  object
17   extraction_type_class                  44550 non-null  object
18   management                             44550 non-null  object
19   management_group                       44550 non-null  object
20   payment                                44550 non-null  object
21   water_quality                          44550 non-null  object
22   quantity                               44550 non-null  object
23   source                                 44550 non-null  object
24   source_class                           44550 non-null  object
25   waterpoint_type_group                  44550 non-null  object
26   month_recorded                         44550 non-null  int64
dtypes: bool(2), float64(3), int64(7), object(15)
memory usage: 8.9+ MB
```

In [427]:

```
1
2
3 # create subpipe for numeric data
4
5 subpipe_num = Pipeline(steps=[('num_impute', SimpleImputer()),
6                               ('ss', StandardScaler())])
7
8 # create subpipe for categorical data, use SimpleImputer for 'miss
9
10 subpipe_cat = Pipeline(steps=[('cat_impute', SimpleImputer(strateg
11                             ('ohe', OneHotEncoder(sparse=False, h
12
13 # combine subpipes into ColumnTransformer
14
15 CT = ColumnTransformer(transformers=[('subpipe_num', subpipe_num,
16                                     ('subpipe_cat', subpipe_cat, [
17
18
19
20                                     remainder='passthrough')
21
22
```

In [428]:

```
1 #Perform Logistic Regression for initial model
2
3 log_reg_pipe = Pipeline(steps = [('ct', CT),
4                                  ('log_reg', LogisticRegression(random_s
```

```
In [429]: 1 # Fit the logistic regression model
          2 log_reg_pipe.fit(X_train, y_train)
```

```
Out[429]: Pipeline(steps=[('ct',
                             ColumnTransformer(remainder='passthrough',
                                                    transformers=[('subpipe_num',
                                                                    Pipeline(steps=[('n
um_impute',
                                                                    Si
mpleImputer()),
                                                                    ('s
s',
                                                                    St
andardScaler())])),
                             [0, 2, 4, 5, 12]),
                            ('subpipe_cat',
                             Pipeline(steps=[('c
at_impute',
                                                                    Si
mpleImputer(fill_value='missing',
                                                                    strategy='constant'))],
                                                                    ('o
he',
                                                                    On
eHotEncoder(handle_unknown='ignore',
                                                                    sparse=False))])),
                             [1, 3, 6, 7, 8, 9,
                             10, 11, 13,
                             14, 15, 16, 17, 1
                             8, 19, 20,
                             21, 22, 23, 24, 2
                             5, 26])])),
                  ('log_reg', LogisticRegression(random_state=42)))]
```

Evaluate initial model

```
In [430]: 1 # Score the log reg model
          2 log_reg_pipe.score(X_train, y_train)
```

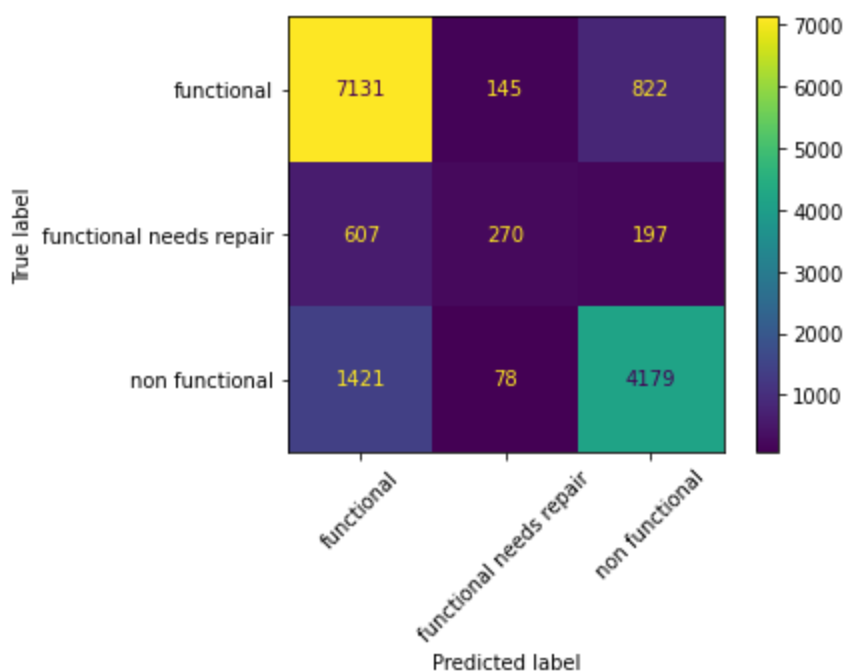
```
Out[430]: 0.8025813692480359
```

```
In [431]: 1 # create predicted target variable
          2 y_hat = log_reg_pipe.predict(X_test)
```

```
In [432]: 1 # Generate log_reg classification report
          2 from sklearn.metrics import classification_report
          3
          4 print(classification_report(y_test, y_hat))
```

	precision	recall	f1-score	support
functional	0.78	0.88	0.83	8098
functional needs repair	0.55	0.25	0.34	1074
non functional	0.80	0.74	0.77	5678
accuracy			0.78	14850
macro avg	0.71	0.62	0.65	14850
weighted avg	0.77	0.78	0.77	14850

```
In [433]: 1 plot_confusion_matrix(log_reg_pipe, X_test, y_test, xticks_rotatio
```



```
In [434]: 1 # Save model in Joblib
2 from joblib import Parallel, delayed
3 import joblib
4
5 import pickle
6
7 # Save the model as a pickle in a file
8 joblib.dump(log_reg_pipe, 'log_reg.pkl')
9
10 # Load the model from the file
11 #log_reg_from_joblib = joblib.load('log_reg.pkl')
12
13 # Use the loaded model to make predictions
14 #log_reg_from_joblib.predict(X_test)
```

```
Out[434]: ['log_reg.pkl']
```

3. Decision Tree Model with Parameter Tuning

Considering the dataset a decision tree would be a useful secondary model. Use hyperparameter tuning to improve upon the initial logistic regression model.

```
1 import category_encoders as ce
2
3
4 # create subpipe for numeric data
5
6 subpipe_num = Pipeline(steps=[('num_impute', SimpleImputer()),
7                               ('ss', StandardScaler())])
8
9 # create subpipe for categorical data, use SimpleImputer for
  'missing' data.
10
11 subpipe_cat = Pipeline(steps=[('cat_impute',
12                               SimpleImputer(strategy='constant', fill_value = 'missing')),
13                               ('ohe', OneHotEncoder(sparse=False,
14                                                       handle_unknown='ignore'))])
14
15 # combine subpipes into ColumnTransformer
16 CT = ColumnTransformer(transformers=[('subpipe_num', subpipe_num,
17 [0, 2, 4, 5, 12]),
18                               ('subpipe_cat', subpipe_cat,
19 [1, 3, 6, 7, 8, 9, 10,
20 11, 13, 14, 15, 16,
21 17, 18, 19, 20, 21,
22 22, 23, 24, 25, 26])],
23                               remainder='passthrough', n_jobs = -1)
24
```

```
In [549]: 1 import category_encoders as ce
2
3
4 # create subpipe for numeric data
5
6 subpipe_num = Pipeline(steps=[('num_impute', SimpleImputer()),
7                               ('ss', StandardScaler())])
8
9 # create subpipe for categorical data, use SimpleImputer for 'miss
10
11 subpipe_cat = Pipeline(steps=[('cat_impute', SimpleImputer(strateg
12                               ('ce_loo', ce.OrdinalEncoder(return_d
13
14 # combine subpipes into ColumnTransformer
15
16 CT_loo = ColumnTransformer(transformers=[('subpipe_num', subpipe_n
17                                         ('subpipe_cat', subpipe_cat, [
18
19
20
21                                         remainder='passthrough', n_jobs = -1)
22
```

```
In [550]: 1 from sklearn.preprocessing import LabelEncoder
2
3 le = LabelEncoder()
4 le.fit_transform(y_train)
```

Out[550]: array([2, 0, 0, ..., 2, 0, 2])

```
In [551]: 1 cols = ['funder', 'installer', 'num_private', 'basin',
2            'region_code', 'district_code', 'lga', 'ward',
3            'construction_year', 'extraction_type', 'extraction_type_cl
4            'management', 'management_group', 'payment', 'water_quality
5            'quantity', 'source', 'source_class', 'waterpoint_type_grou
6            'month_recorded']
7
8 X_train[cols] = X_train[cols].astype(str)
9 X_test[cols] = X_test[cols].astype(str)
```

```
In [552]: 1 # Use a decision tree for the secondary model
2 dtc = DecisionTreeClassifier(random_state = 42)
3
4 dtc_pipe = Pipeline(steps=[('ct_loo', CT_loo),
5                             ('dtc', dtc)])
```

```
In [553]: 1 dtc_pipe.fit(X_train, y_train)
```

```
Out[553]: Pipeline(steps=[('ct_loo',
                             ColumnTransformer(n_jobs=-1, remainder='passthrough',
                                                    transformers=[('subpipe_num',
                                                                    Pipeline(steps=[('n
um_impute',
                                                                    Si
mpleImputer()),
                                                                    ('s
s',
                                                                    St
andardScaler())])),
                             [0, 2, 4, 5, 12]),
                             ('subpipe_cat',
                              Pipeline(steps=[('c
at_impute',
                              Si
mpleImputer(fill_value='missing',
strategy='constant'))),
                              ('c
e_loo',
                              Or
dinalEncoder())])),
                             [1, 3, 6, 7, 8, 9,
                             10, 11, 13,
                             14, 15, 16, 17, 1
                             8, 19, 20,
                             21, 22, 23, 24, 2
                             5, 26])])),
              ('dtc', DecisionTreeClassifier(random_state=42)))]
```

```
In [554]: 1 dtc_pipe.score(X_train, y_train)
```

```
Out[554]: 0.9984960718294051
```

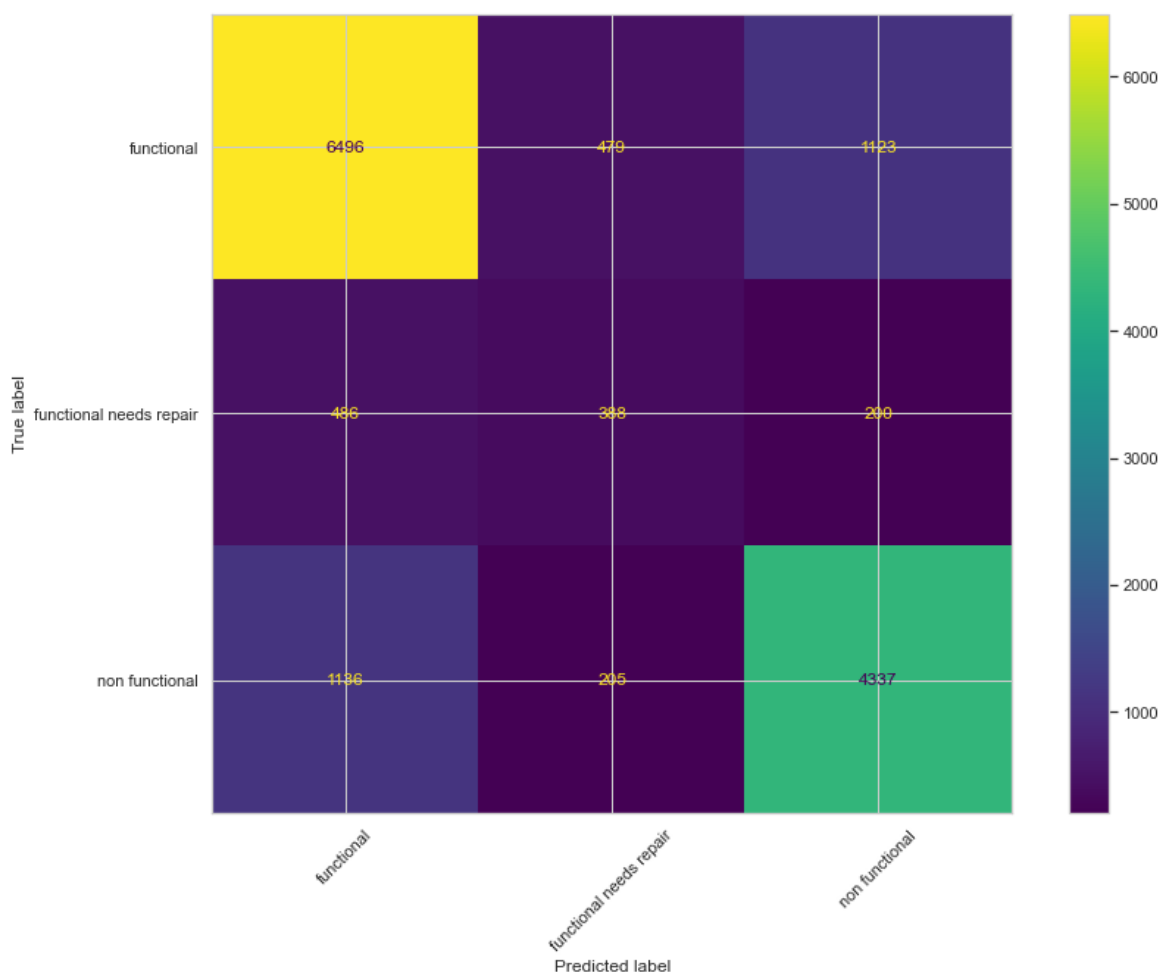
Evaluate Decision Tree Model

```
In [555]: 1 y_hat = dtc_pipe.predict(X_test)
```


In [556]: 1 `print(classification_report(y_test, y_hat))`

	precision	recall	f1-score	support
functional	0.80	0.80	0.80	8098
functional needs repair	0.36	0.36	0.36	1074
non functional	0.77	0.76	0.77	5678
accuracy			0.76	14850
macro avg	0.64	0.64	0.64	14850
weighted avg	0.76	0.76	0.76	14850

In [557]: 1 `plot_confusion_matrix(dtc_pipe, X_test, y_test, xticks_rotation=45)`



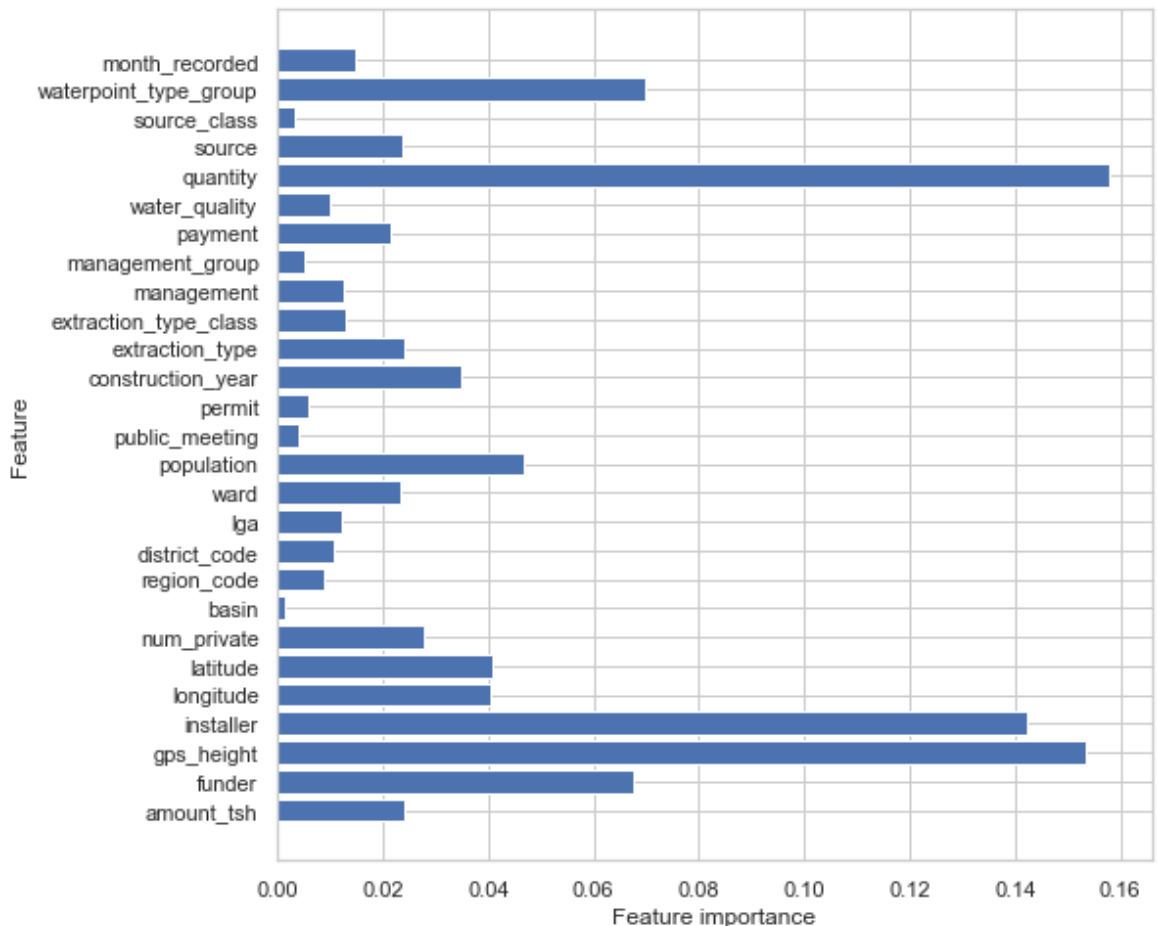
In [558]: 1 `len(dtc_pipe.named_steps['dtc'].feature_importances_)`

Out [558]: 27

```
In [559]: 1 model_tree = dtc_pipe.named_steps['dtc']
          2 model_tree.feature_importances_
```

```
Out[559]: array([0.02420042, 0.06744788, 0.15358044, 0.14229542, 0.04049446,
                  0.04084094, 0.02766165, 0.0013189 , 0.00897386, 0.01082605,
                  0.01215422, 0.02326186, 0.04666773, 0.00409639, 0.00569057,
                  0.0350141 , 0.02393515, 0.01301596, 0.01250364, 0.00512469,
                  0.02150802, 0.00995001, 0.15799101, 0.02366923, 0.00320772,
                  0.06980831, 0.01476137])
```

```
In [560]: 1 def plot_feature_importances(model):
          2     n_features = X_train.shape[1]
          3     plt.figure(figsize=(8,8))
          4     plt.barh(range(n_features), model.feature_importances_, align=
          5     plt.yticks(np.arange(n_features), X_train.columns.values)
          6     plt.xlabel('Feature importance')
          7     plt.ylabel('Feature')
          8
          9     plot_feature_importances(model_tree)
```



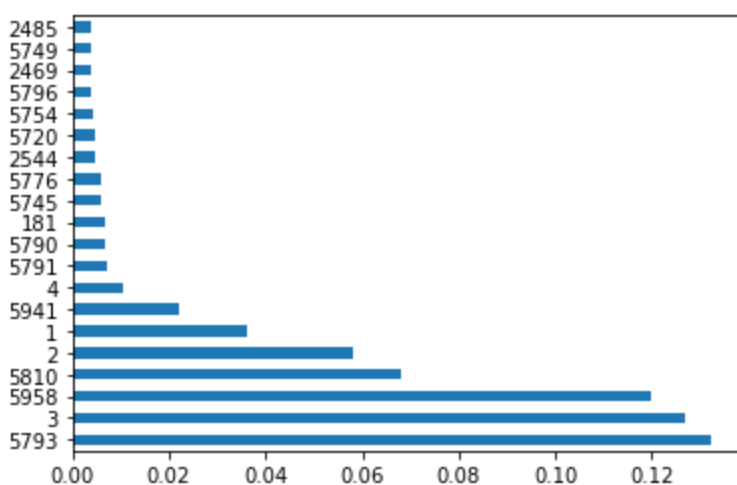
```
In [ ]: 1 # Conduct inverse_transform of one hot encoding
```

```
In [528]: 1 important_features_dict = {}
          2 for idx, val in enumerate(model_tree.feature_importances_):
          3     important_features_dict[idx] = val
          4
          5 important_features_list = sorted(important_features_dict,
          6                                 key=important_features_dict.get,
          7                                 reverse=True)
          8
          9 print(f'5 most important features: {important_features_list[:10]}')
```

5 most important features: [5941, 2, 3, 5958, 1, 4, 0, 5, 5939, 5938]

```
In [540]: 1 feat_importances = pd.Series(model_tree.feature_importances_, index=important_features_list)
          2 feat_importances.nlargest(20).plot(kind='barh')
```

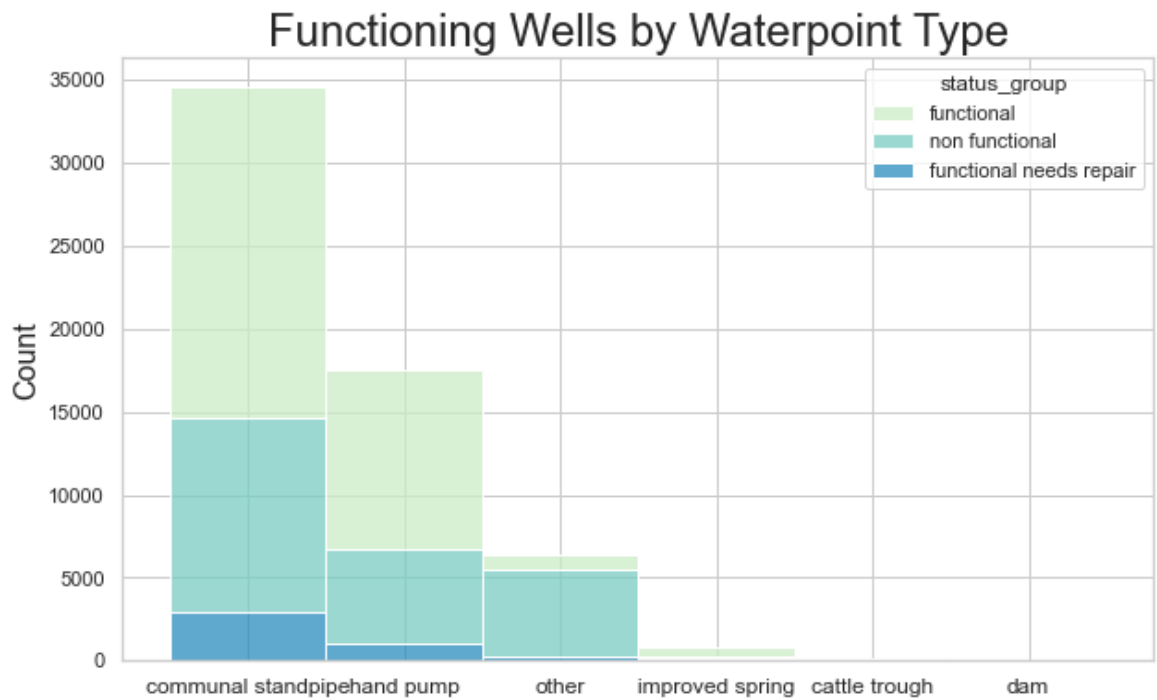
Out[540]: <AxesSubplot:>



```
In [530]: 1 X_train.nunique()
```

```
Out[530]: amount_tsh          95
          funder             1645
          gps_height          2391
          installer          1861
          longitude          43171
          latitude           43173
          num_private         59
          basin               9
          region_code         27
          district_code       20
          lga                 125
          ward                2071
          population          951
          public_meeting       2
          permit              2
          construction_year    55
          extraction_type      18
          extraction_type_class 7
          management          12
          management_group     5
          payment              7
          water_quality         8
          quantity             5
          source              10
          source_class         3
          waterpoint_type_group 6
          month_recorded       12
          dtype: int64
```

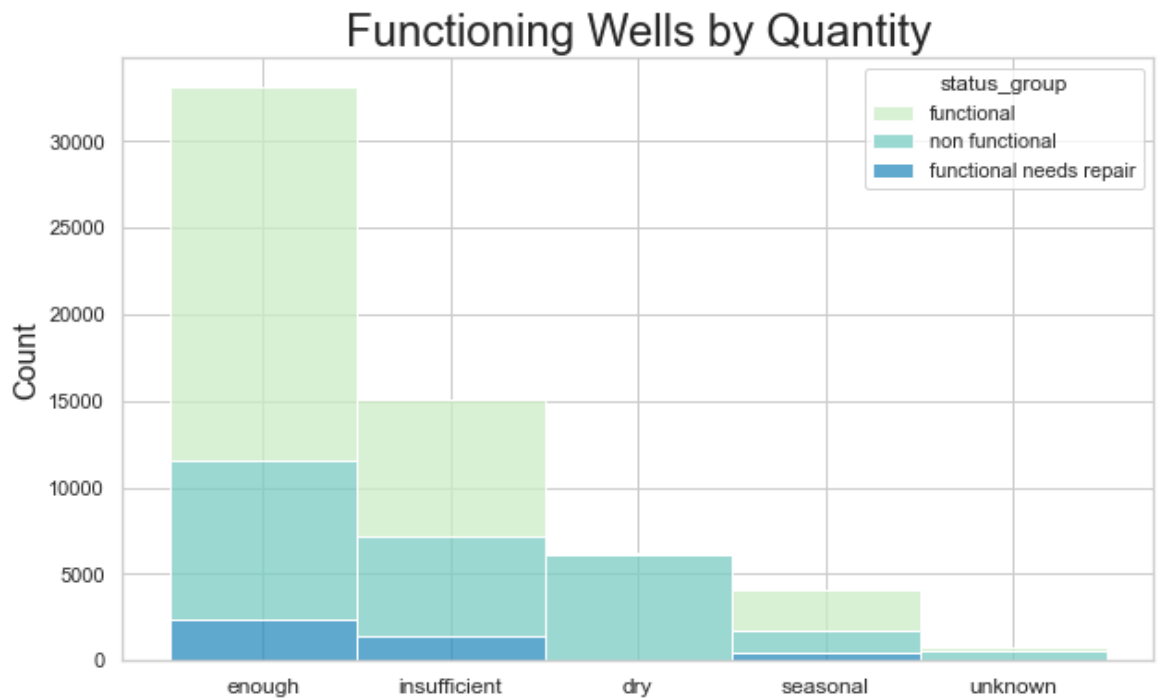
```
In [543]: 1 sns.set_theme()
2
3 sns.set(rc={"figure.figsize":(10, 6)})
4 sns.set_style('whitegrid')
5
6 sns.histplot(data = df, x = 'waterpoint_type_group', hue = 'status
7             bins = 10, binwidth = 6, palette = 'GnBu', legend = '
8             multiple = 'stack')
9
10
11 plt.title("Functioning Wells by Waterpoint Type", fontsize= 24)
12 plt.xlabel(None)
13 plt.ylabel("Count", fontsize = 16)
14 plt.xticks(rotation = 0, fontsize = 12);
```



```

In [545]: 1 sns.set_theme()
          2
          3 sns.set(rc={"figure.figsize":(10, 6)})
          4 sns.set_style('whitegrid')
          5
          6 sns.histplot(data = df, x = 'quantity', hue = 'status_group',
          7               bins = 10, binwidth = 6, palette = 'GnBu', legend = '
          8               multiple = 'stack')
          9
         10
         11 plt.title("Functioning Wells by Quantity",fontsize= 24)
         12 plt.xlabel(None)
         13 plt.ylabel("Count", fontsize = 16)
         14 plt.xticks(rotation = 0, fontsize = 12);

```



```

In [558]: 1 # Save the decision tree model as a pickle in a file
          2 joblib.dump(dtc_pipe, 'dtc_pipe.pkl')

```

Out[558]: ['dtc_pipe.pkl']

Results

The decision tree model's accuracy was less than the logistic regression model and did not improve upon the logistic regression accuracy, though the f1-score for non-functional wells of both models remained the same.

Use gridsearch for hyperparameter tuning.

```
In [75]: 1 params = {}
          2 params['dtc__criterion'] = ['gini', 'entropy']
          3 params['dtc__min_samples_leaf'] = [1, 3, 5, 7, 10]
          4 params['dtc__max_depth'] = [1,3,5,7,9]
          5 params['dtc__splitter'] = ['best', 'random']
          6
          7 gs = GridSearchCV(estimator=dtc_pipe,
          8                     param_grid=params,
          9                     cv=3)
```

```
In [76]: 1 gs.fit(X_train, y_train)
```

```
Out[76]: GridSearchCV(cv=3,
                      estimator=Pipeline(steps=[('ct',
                                                  ColumnTransformer(remainder='
passthrough',
                                                                transformer
s=[('subpipe_num',
Pipeline(steps=[('num_impute',
SimpleImputer()),
('ss',
StandardScaler())])),
[0, 2,
4, 5,
12]),
('subpipe_cat',
Pipeline(steps=[('cat_impute',
SimpleImputer(fill_value='missing',
strategy='constant'))),
('ohe',
OneHotEncoder(handle_unknown='ignore',
sparse=False))])),
[1, 3,
6, 7,
8, 9,
10,
11,
13,
14,
15,
16,
17,
```



```
18,  
19,  
20,  
21,  
22,  
23,  
24,  
25,  
26]]))),  
                                ('dtc',  
                                DecisionTreeClassifier(random  
_state=42))]),  
                                param_grid={'dtc__criterion': ['gini', 'entropy'],  
                                            'dtc__max_depth': [1, 3, 5, 7, 9],  
                                            'dtc__min_samples_leaf': [1, 3, 5, 7, 10],  
                                            'dtc__splitter': ['best', 'random']})
```

```
In [77]: 1 # Identify the best parameters  
        2 gs.best_params_
```

```
Out[77]: {'dtc__criterion': 'gini',  
          'dtc__max_depth': 9,  
          'dtc__min_samples_leaf': 5,  
          'dtc__splitter': 'best'}
```

```
In [78]: 1 # Examine cross validation results
          2 gs.cv_results_['mean_test_score']
```

```
Out[78]: array([0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 ,
                 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 ,
                 0.69427609, 0.69441077, 0.69427609, 0.69441077, 0.69429854,
                 0.69443322, 0.69427609, 0.69445567, 0.69450056, 0.69461279,
                 0.71270483, 0.70826038, 0.71261504, 0.70810325, 0.71265993,
                 0.70808081, 0.71272727, 0.7081257 , 0.71261504, 0.70808081,
                 0.72489338, 0.72282828, 0.72480359, 0.72255892, 0.72455668,
                 0.72253648, 0.72453423, 0.72210999, 0.72430976, 0.72190797,
                 0.73719416, 0.73476992, 0.73705948, 0.73429854, 0.73748597,
                 0.73297419, 0.73643098, 0.73313131, 0.73542088, 0.73236813,
                 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 ,
                 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 , 0.6424018 ,
                 0.69342312, 0.6935578 , 0.69342312, 0.6935578 , 0.69342312,
                 0.6935578 , 0.69342312, 0.69360269, 0.69351291, 0.69362514,
                 0.7006734 , 0.69723906, 0.70060606, 0.6973064 , 0.70056117,
                 0.6973064 , 0.70042649, 0.69717172, 0.70038159, 0.69710438,
                 0.71353535, 0.71380471, 0.71353535, 0.71360269, 0.71384961,
                 0.71378227, 0.71367003, 0.7138945 , 0.71335578, 0.71411897,
                 0.73158249, 0.7308642 , 0.73167228, 0.7308642 , 0.73156004,
                 0.7308193 , 0.73182941, 0.73021324, 0.73113356, 0.73005612])
```

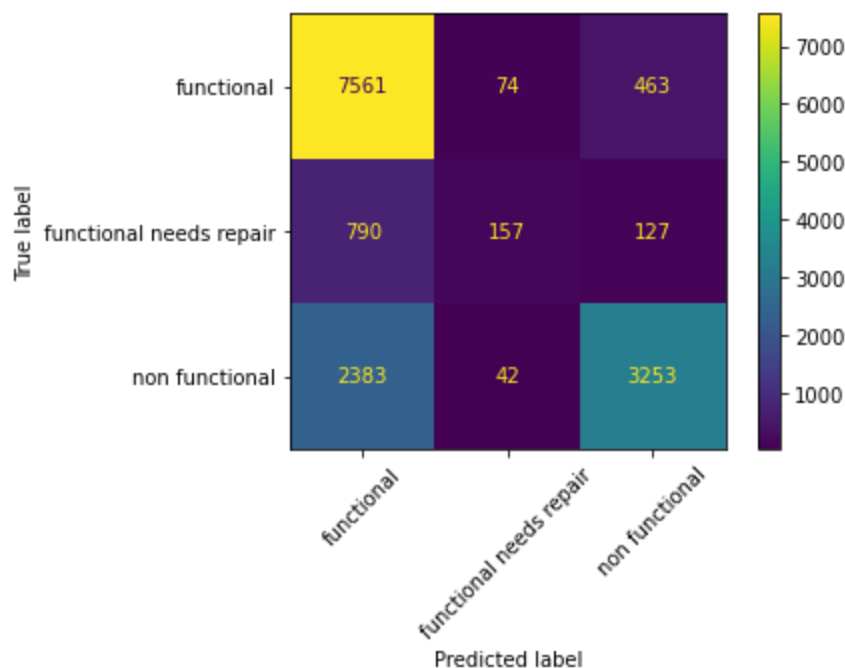
Evaluate decision tree gridsearch results

```
In [79]: 1 y_hat = gs.predict(X_test)
```

```
In [80]: 1 print(classification_report(y_test, y_hat))
```

	precision	recall	f1-score	support
functional	0.70	0.93	0.80	8098
functional needs repair	0.58	0.15	0.23	1074
non functional	0.85	0.57	0.68	5678
accuracy			0.74	14850
macro avg	0.71	0.55	0.57	14850
weighted avg	0.75	0.74	0.72	14850

```
In [81]: 1 plot_confusion_matrix(gs, X_test, y_test, xticks_rotation = 45);
```



```
In [559]: 1 # Save the model as a pickle in a file  
2 joblib.dump(gs, 'grid_search_dtc.pkl')
```

```
Out[559]: ['grid_search_dtc.pkl']
```

Results

While accuracy decreased overall, the precision score on non-functional wells improved from 77% to 85%. This could be a good model if we only focus on precision score for non-functioning wells. Wells that need repair precision score also improved by 20%, this opens a path to possibly identify wells that could soon be non-functioning.

4. Random Forest with SMOTE and Tuning

Use a random forest model to further explore whether the precision or recall score on non-functioning wells can be improved. Address class imbalance issues with SMOTE. Further tune the model using search tools for best hyperparameters.

```
In [562]: 1 # Instantiate a Random Forest Classifier
          2
          3 rfc = RandomForestClassifier(random_state=42)
          4
          5 # Instantiate SMOTE for class imbalance
          6
          7 sm = SMOTE(sampling_strategy = 'auto', random_state = 42)
          8
          9 # Create pipeline
         10
         11 rfc_model_pipe = ImPipeline(steps=[('ct_loo', CT_loo),
         12                                   ('sm', sm),
         13                                   ('rfc', rfc)])
         14
```

```
In [563]: 1 rfc_model_pipe.fit(X_train, y_train)
```

```
Out[563]: Pipeline(steps=[('ct_loo',
                             ColumnTransformer(n_jobs=-1, remainder='passthrough',
                                                  transformers=[('subpipe_num',
                                                                Pipeline(steps=[('num_impute',
                                                                                          SimpleImputer()),
                                                                                          ('standardScaler', StandardScaler())])),
                                                                [0, 2, 4, 5, 12]),
                                                                ('subpipe_cat',
                                                                Pipeline(steps=[('cat_impute',
                                                                                          SimpleImputer(fill_value='missing',
                                                                                          strategy='constant'))],
                                                                ('ct_loo',
                                                                OrdinalEncoder()))],
                                                                [1, 3, 6, 7, 8, 9,
                                                                10, 11, 13,
                                                                14, 15, 16, 17, 18, 19, 20,
                                                                21, 22, 23, 24, 25, 26])])),
                             ('sm', SMOTE(random_state=42)),
                             ('rfc', RandomForestClassifier(random_state=42))])
```

```
In [564]: 1 rfc_model_pipe.score(X_train, y_train)
```

```
Out[564]: 0.9984960718294051
```

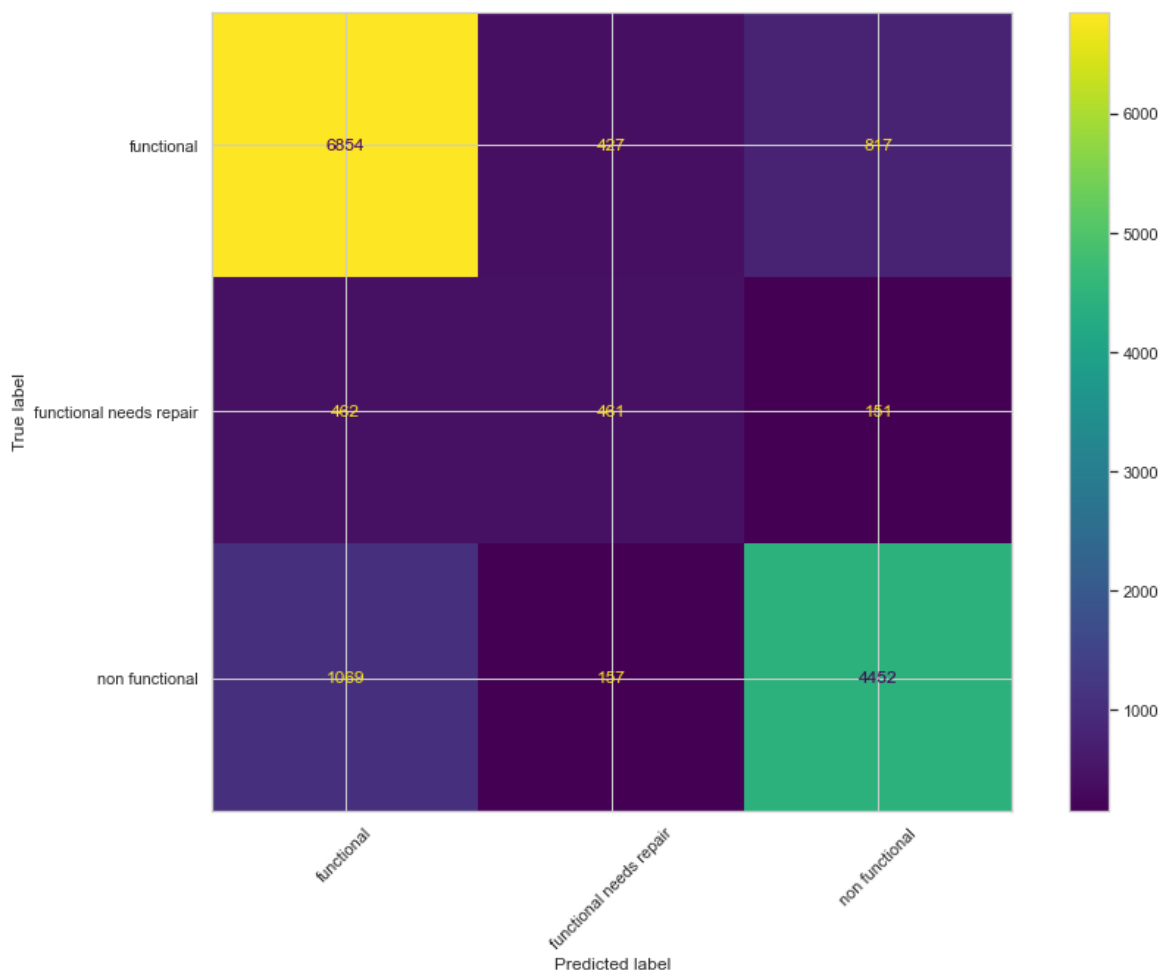
Evaluate results on Random Forest

```
In [565]: 1 y_hat_rfc = rfc_model_pipe.predict(X_test)
```

```
In [566]: 1 print(classification_report(y_test, y_hat_rfc))
```

	precision	recall	f1-score	support
functional	0.82	0.85	0.83	8098
functional needs repair	0.44	0.43	0.44	1074
non functional	0.82	0.78	0.80	5678
accuracy			0.79	14850
macro avg	0.69	0.69	0.69	14850
weighted avg	0.79	0.79	0.79	14850

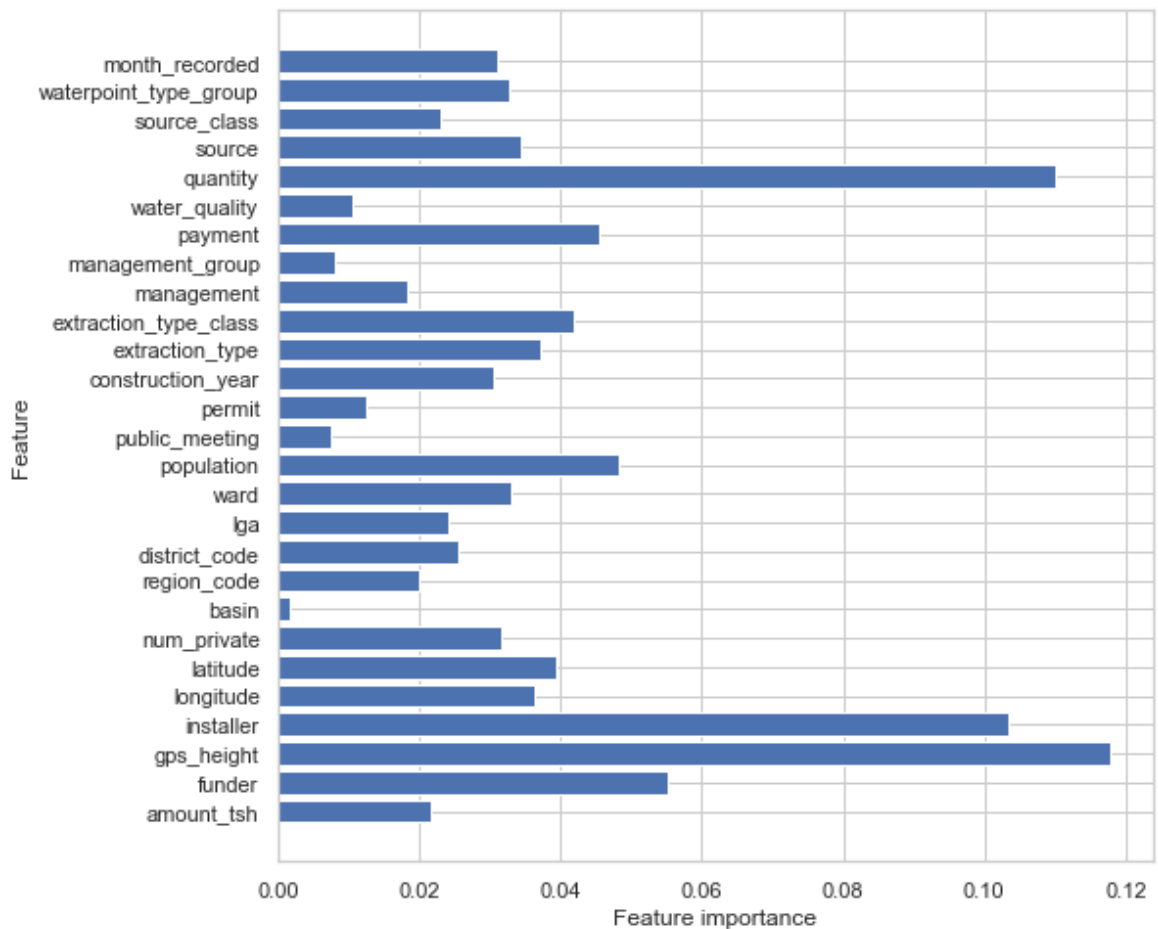
```
In [567]: 1 plot_confusion_matrix(rfc_model_pipe, X_test, y_test, xticks_rotat
```



```
In [568]: 1 # Save the random forest model as a pickle in a file
          2 joblib.dump(rfc_model_pipe, 'rfc_model.pkl')
```

```
Out[568]: ['rfc_model.pkl']
```

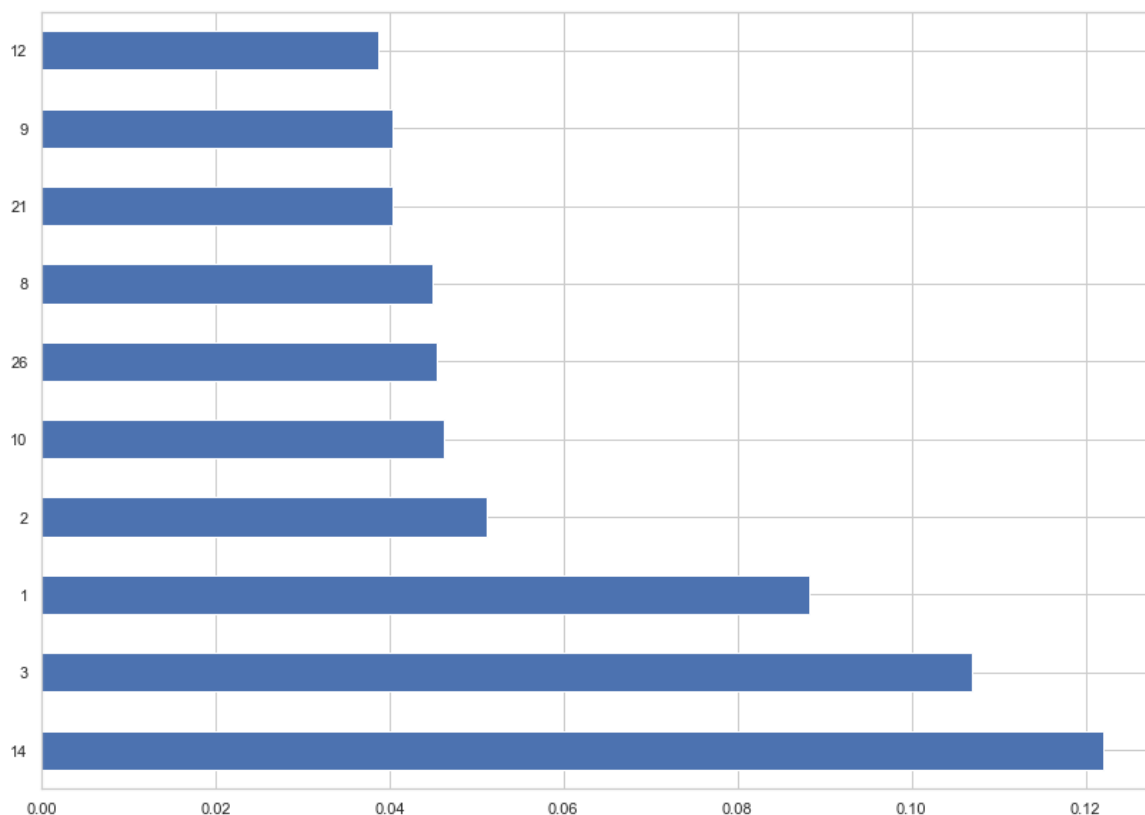
```
In [587]: 1 def plot_feature_importances(model):
2         n_features = X_train.shape[1]
3         plt.figure(figsize=(8,8))
4         plt.barh(range(n_features), model.feature_importances_, align=
5         plt.yticks(np.arange(n_features), X_train.columns.values)
6         plt.xlabel('Feature importance')
7         plt.ylabel('Feature')
8
9         plot_feature_importances(rfc_model_pipe.named_steps['rfc'])
10
11
```



```
In [579]: 1 important_features_dict = {}
2         for idx, val in enumerate(final_model.named_steps['rfc'].feature_i
3             important_features_dict[idx] = val
4
5         important_features_list = sorted(important_features_dict,
6                                         key=important_features_dict.get,
7                                         reverse=True)
8
9         print(f'5 most important features: {important_features_list[:10]}')
5 most important features: [22, 2, 3, 1, 17, 12, 20, 23, 16, 5]
```

```
In [582]: 1 feat_importances = pd.Series(final_model.named_steps['rfc'].feature
2         feat_importances.nlargest(10).plot(kind='barh')
```

Out[582]: <AxesSubplot:>



Gridsearch for hyperparameter tuning

```
In [95]: 1 # Grid Search for better model criteria
2
3 params = {'rfc__n_estimators': [10],
4           'rfc__criterion': ['gini'],
5           'rfc__min_samples_leaf': [1, 5, 10],
6           'rfc__max_depth': [1, 5, 9],
7           'rfc__max_features': [9]
8           }
9
10 gs_rfc = GridSearchCV(estimator=rfc_model_pipe,
11                       param_grid=params, n_jobs = -1,
12                       cv=3)
```

```
In [96]: 1 gs_rfc.fit(X_train, y_train)
```

```
Out[96]: GridSearchCV(cv=3,
                      estimator=Pipeline(steps=[('ct',
                                                  ColumnTransformer(remainder='
passthrough',
                                                                transformer
s=[('subpipe_num',
Pipeline(steps=[('num_impute',
SimpleImputer()),
('ss',
StandardScaler())])),
[0, 2,
4, 5,
12]),
('subpipe_cat',
Pipeline(steps=[('cat_impute',
SimpleImputer(fill_value='missing',
strategy='constant'))),
('ohe',
OneHotEncoder(handle_unknown='ignore',
sparse=False))])),
[1, 3,
6, 7,
8, 9,
10,
11,
13,
14,
15,
16,
17,
```



```

18,
19,
20,
21,
22,
23,
24,
25,
26]]))),
                                ('sm', SMOTE(random_state=4
2)),
                                ('rfc',
                                RandomForestClassifier(random
_state=42))]),
                                n_jobs=-1,
                                param_grid={'rfc__criterion': ['gini'],
                                'rfc__max_depth': [1, 5, 9], 'rfc__max_featu
res': [9],
                                'rfc__min_samples_leaf': [1, 5, 10],
                                'rfc__n_estimators': [10]})

```

```

In [103]: 1 # Best parameters for further tuning
          2 gs_rfc.best_params_

```

```

Out[103]: {'rfc__criterion': 'gini',
           'rfc__max_depth': 9,
           'rfc__max_features': 9,
           'rfc__min_samples_leaf': 5,
           'rfc__n_estimators': 10}

```

```

In [104]: 1 gs_rfc.score(X_train, y_train)

```

```

Out[104]: 0.5083726150392817

```

Evaluate gridsearch results

```

In [105]: 1 y_hat_gs_rfc = gs_rfc.predict(X_test)

```

In [106]: 1 `print(classification_report(y_test, y_hat_gs_rfc))`

	precision	recall	f1-score	support
functional	0.70	0.50	0.58	8098
functional needs repair	0.13	0.56	0.21	1074
non functional	0.62	0.48	0.54	5678
accuracy			0.50	14850
macro avg	0.48	0.51	0.44	14850
weighted avg	0.63	0.50	0.54	14850

In [561]: 1 `# Save the rfc gridsearch model model as a pickle in a file`
 2 `joblib.dump(gs_rfc, 'gs_rfc.pkl')`

Out[561]: ['gs_rfc.pkl']

Results summary on Random Forest Gridsearch

Accuracy decreased significantly, perhaps as a result of using 10 `n_estimators` rather than the default 100 to cut down on processing time. This model though suggests where to explore for `max_depth`, and `samples leaf` and `split`.

Use Randomized Search

```
In [137]: 1
2 from sklearn.model_selection import RandomizedSearchCV
3
4 # Based in previous gridsearch, optimize for max depth, min sample
5
6 random_grid = {
7     'rfc__bootstrap': [True],
8     'rfc__max_depth': [10, 20, 50, 100],
9     'rfc__max_features': ['auto', 'sqrt'],
10    'rfc__min_samples_leaf': [1, 2, 4],
11    'rfc__min_samples_split': [2, 5, 10],
12    'rfc__n_estimators': [10, 100]
13 }
14
15 random = RandomizedSearchCV(estimator = rfc_model_pipe,
16                             param_distributions = random_grid, n_jobs = -1,
17                             verbose = 2, random_state = 42, cv=3)
```

In [138]: 1 random.fit(X_train, y_train)

Fitting 3 folds for each of 10 candidates, totalling 30 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 124.2min finished

Out[138]: RandomizedSearchCV(cv=3,
 estimator=Pipeline(steps=[('ct',
 ColumnTransformer(remainders='passthrough',
 transformers=[('subpipe_num',
 Pipeline(steps=[('num_impute',
 SimpleImputer()),
 ('ss',
 StandardScaler())]),
 [0,
 2,
 4,
 5,
 12]),
 ('subpipe_cat',
 Pipeline(steps=[('cat_impute',
 SimpleImputer(fill_value='missing',
 strategy='constant'))],
 ('ohe',
 OneHotEncoder(handle_unknown='drop',
 20,
 21,
 22,
 23,
 24,

```

25,
26]]]]),
                                ('sm', SMOTE(random_state=42)),
                                ('rfc',
                                 RandomForestClassifier
                                 (random_state=42))),
                                n_jobs=-1,
                                param_distributions={'rfc__bootstrap': [True],
                                                    'rfc__max_depth': [10, 20, 5
0, 100],
                                                    'rfc__max_features': ['auto',
'sqrt'],
                                                    'rfc__min_samples_leaf': [1,
2, 4],
                                                    'rfc__min_samples_split': [2,
5, 10],
                                                    'rfc__n_estimators': [10, 10
0]}},
                                random_state=42, verbose=2)

```

```

In [140]: 1 # Best paramters from randomized search on RFC
          2 random.best_params_

```

```

Out[140]: {'rfc__n_estimators': 100,
           'rfc__min_samples_split': 5,
           'rfc__min_samples_leaf': 2,
           'rfc__max_features': 'auto',
           'rfc__max_depth': 100,
           'rfc__bootstrap': True}

```

Evaluate randomized search best parameter results

```

In [139]: 1 random.score(X_train, y_train)

```

```

Out[139]: 0.8312457912457912

```

```

In [141]: 1 y_hat_random = random.predict(X_test)

```

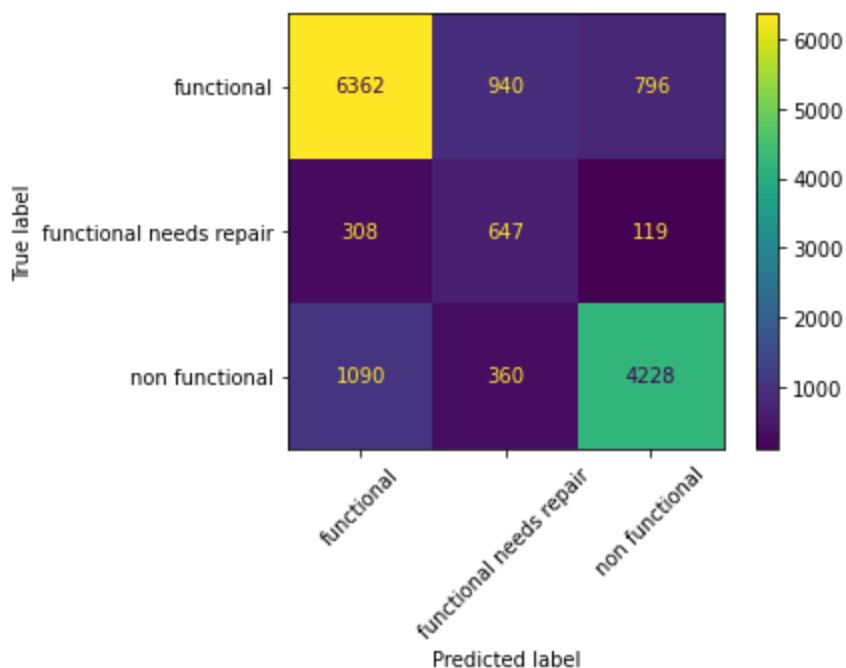
```

In [142]: 1 print(classification_report(y_test, y_hat_random))

```

	precision	recall	f1-score	support
functional	0.82	0.79	0.80	8098
functional needs repair	0.33	0.60	0.43	1074
non functional	0.82	0.74	0.78	5678
accuracy			0.76	14850
macro avg	0.66	0.71	0.67	14850
weighted avg	0.79	0.76	0.77	14850

```
In [569]: 1 plot_confusion_matrix(random, X_test, y_test, xticks_rotation = 45
```



```
In [562]: 1 # Save the model as a pickle in a file
          2 joblib.dump(random, 'random_rfc.pkl')
```

```
Out[562]: ['random_rfc.pkl']
```

Gridsearch based on randomized results

```
In [571]: 1 # Based on randomized search conduct one more gridsearch
          2 params = {
          3     'rfc__n_estimators': [100],
          4     'rfc__min_samples_leaf': [2, 3],
          5     'rfc__max_depth': [100, 150],
          6     'rfc__min_samples_split': [3, 5, 7],
          7     'rfc__max_features': ['auto']
          8 }
          9
         10 gs_rfc_2 = GridSearchCV(estimator=rfc_model_pipe,
         11                          param_grid=params, n_jobs = -1,
         12                          verbose = 2, cv = 3)
         13
```

In [572]: 1 gs_rfc_2.fit(X_train, y_train)

Fitting 3 folds for each of 12 candidates, totalling 36 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 36 out of 36 | elapsed: 129.6min finished

```
Out[572]: GridSearchCV(cv=3,
                      estimator=Pipeline(steps=[('ct',
                                                  ColumnTransformer(remainder='
passthrough',
                                                  transformer
s=[('subpipe_num',
Pipeline(steps=[('num_impute',
SimpleImputer()),
('ss',
StandardScaler())]),
[0, 2,
4, 5,
12]),
('subpipe_cat',
Pipeline(steps=[('cat_impute',
SimpleImputer(fill_value='missing',
strategy='constant'))),
('ohe',
OneHotEncoder(handle_unknown='ign...
22,
23,
24,
25,
26])])),
                      ('sm',
                       SMOTE(n_jobs=-1, random_state
=42)),
                      ('rfc',
                       RandomForestClassifier(max_de
```

```

pth=100,
mples_leaf=2,
mples_split=3,
=-1,
_state=42)))),
        n_jobs=-1,
        param_grid={'rfc__max_depth': [100, 150],
                    'rfc__max_features': ['auto'],
                    'rfc__min_samples_leaf': [2, 3],
                    'rfc__min_samples_split': [3, 5, 7],
                    'rfc__n_estimators': [100]},
        verbose=2)

```

```

In [573]: 1 # Score the model on training data
          2 gs_rfc_2.score(X_train, y_train)

```

Out[573]: 0.8340291806958474

```

In [574]: 1 # examine best paramters
          2 gs_rfc_2.best_params_

```

Out[574]: {'rfc__max_depth': 100,
 'rfc__max_features': 'auto',
 'rfc__min_samples_leaf': 2,
 'rfc__min_samples_split': 3,
 'rfc__n_estimators': 100}

```

In [575]: 1 # create predicted target using test set
          2 y_hat_rfc_2 = gs_rfc_2.predict(X_test)

```

```

In [576]: 1 print(classification_report(y_test, y_hat_rfc_2))

```

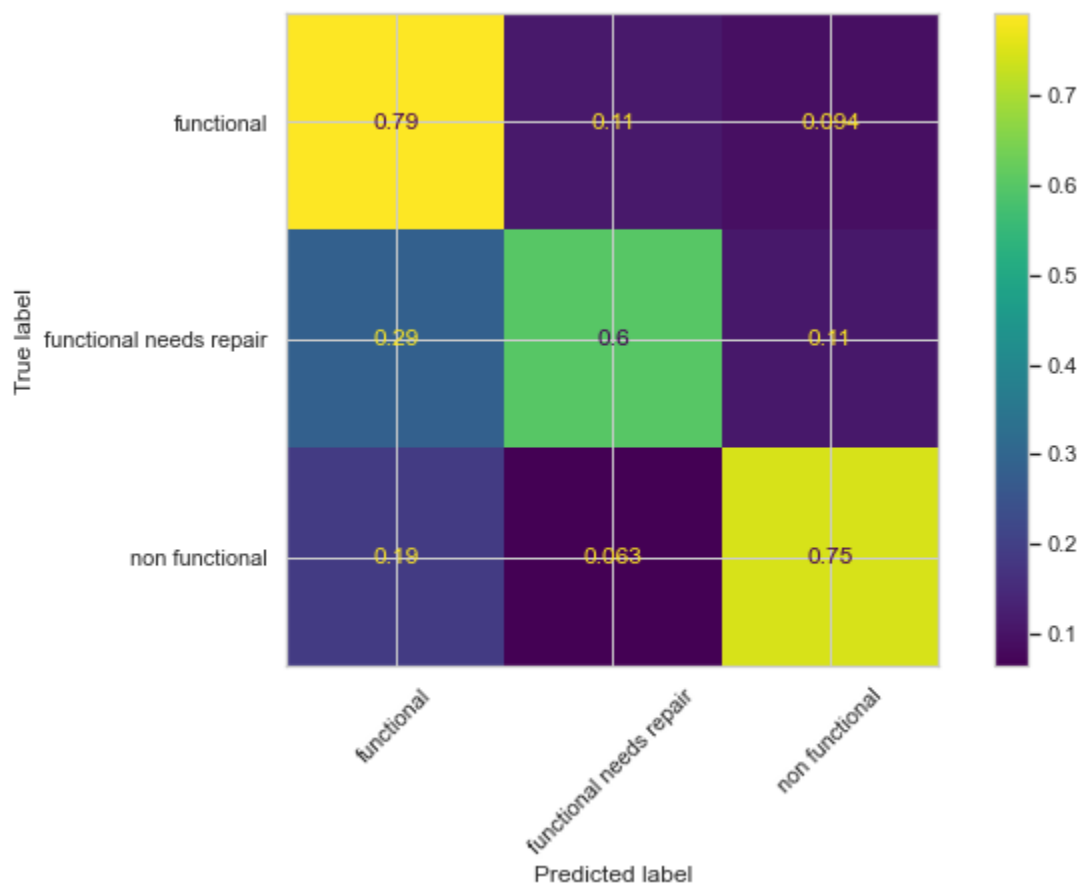
	precision	recall	f1-score	support
functional	0.82	0.79	0.81	8098
functional needs repair	0.34	0.60	0.43	1074
non functional	0.83	0.75	0.79	5678
accuracy			0.76	14850
macro avg	0.66	0.71	0.67	14850
weighted avg	0.79	0.76	0.77	14850

```

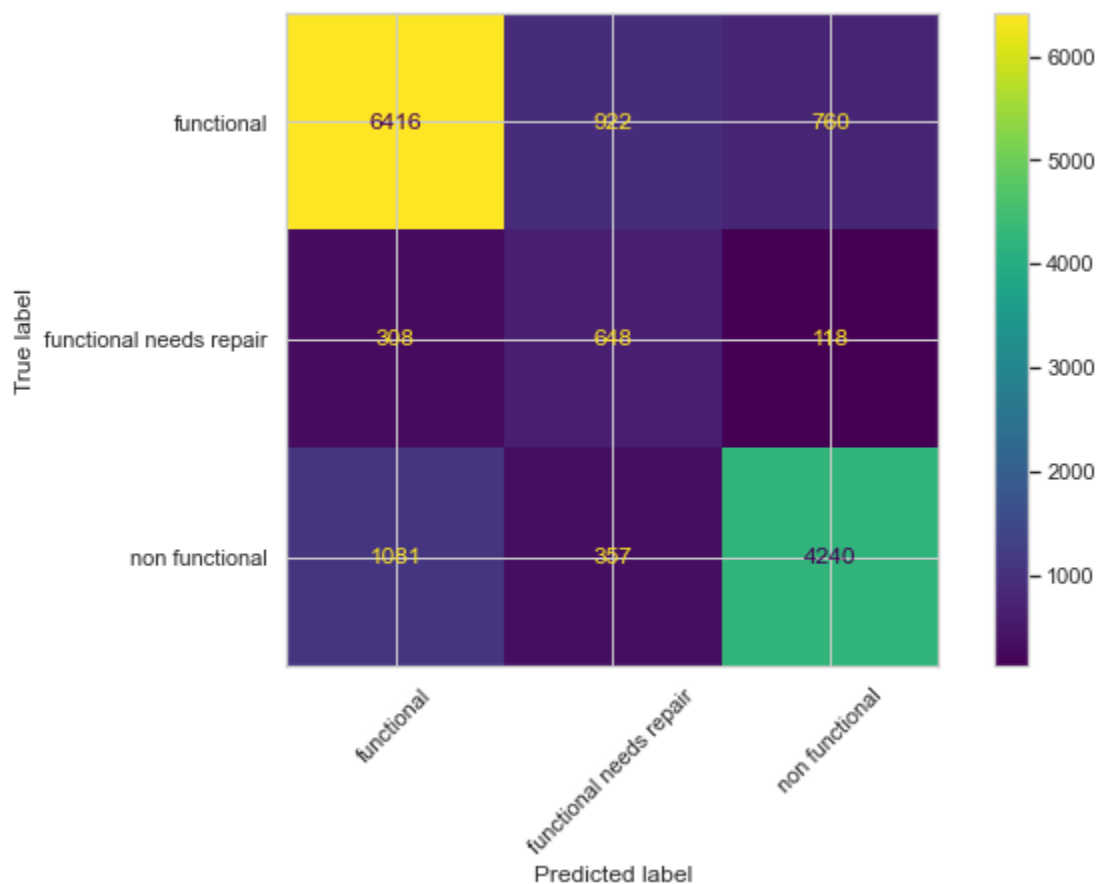
In [577]: 1 # select as final model
          2 final_rfc_model = gs_rfc_2

```

```
In [787]: 1 # Plot confusion matrix with percentages
2 plot_confusion_matrix(final_rfc_model,
3                       X_test, y_test,
4                       xticks_rotation = 45,
5                       normalize = 'true'
6                       );
```




```
In [792]: 1 # Plot matrix with case numbers
2 plot_confusion_matrix(final_rfc_model,
3                       X_test, y_test,
4                       xticks_rotation = 45,
5
6                       );
7
8 plt.savefig('final RFC model matrix')
```



```
In [579]: 1 # Save the model as a pickle in a file
2 joblib.dump(final_rfc_model, 'final_rfc_model.pkl')
```

```
Out[579]: ['final_rfc_model.pkl']
```

```

In [570]: 1 # Instantiate a Random Forest Classifier with final hyperparameter
          2
          3 rfc = RandomForestClassifier(max_depth = 100, max_features = 'auto
          4                                     min_samples_leaf = 2, min_samples_spl
          5                                     random_state=42)
          6
          7 # Instantiate SMOTE for class imbalance
          8
          9 sm = SMOTE(sampling_strategy = 'auto', random_state = 42)
         10
         11 # Create pipeline
         12
         13 final_model = ImPipeline(steps=[('ct_loo', CT_loo),
         14                                   ('sm', sm),
         15                                   ('rfc', rfc)])
         16

```

```

In [571]: 1 final_model.fit(X_train, y_train)

```

```

Out[571]: Pipeline(steps=[('ct_loo',
                             ColumnTransformer(n_jobs=-1, remainder='passthrough',
                             transformers=[('subpipe_num',
                                             Pipeline(steps=[('num_impute',
                                                                    SimpleImputer(),
                                                                    StandardScaler())]),
                                             [0, 2, 4, 5, 12]),
                                             ('subpipe_cat',
                                              Pipeline(steps=[('cat_impute',
                                                                    SimpleImputer(fill_value='missing',
                                                                    strategy='constant'))],
                                             ('c
                                             e_loo',
                                             OrdinalEncoder())]),
                                             [1, 3, 6, 7, 8, 9,
                                             10, 11, 13,
                                             14, 15, 16, 17, 1
                                             8, 19, 20,
                                             21, 22, 23, 24, 2
                                             5, 26])])),
                             ('sm', SMOTE(random_state=42)),
                             ('rfc',
                              RandomForestClassifier(max_depth=100, min_samples_le
                              af=2,
                                                         random_state=42))])

```

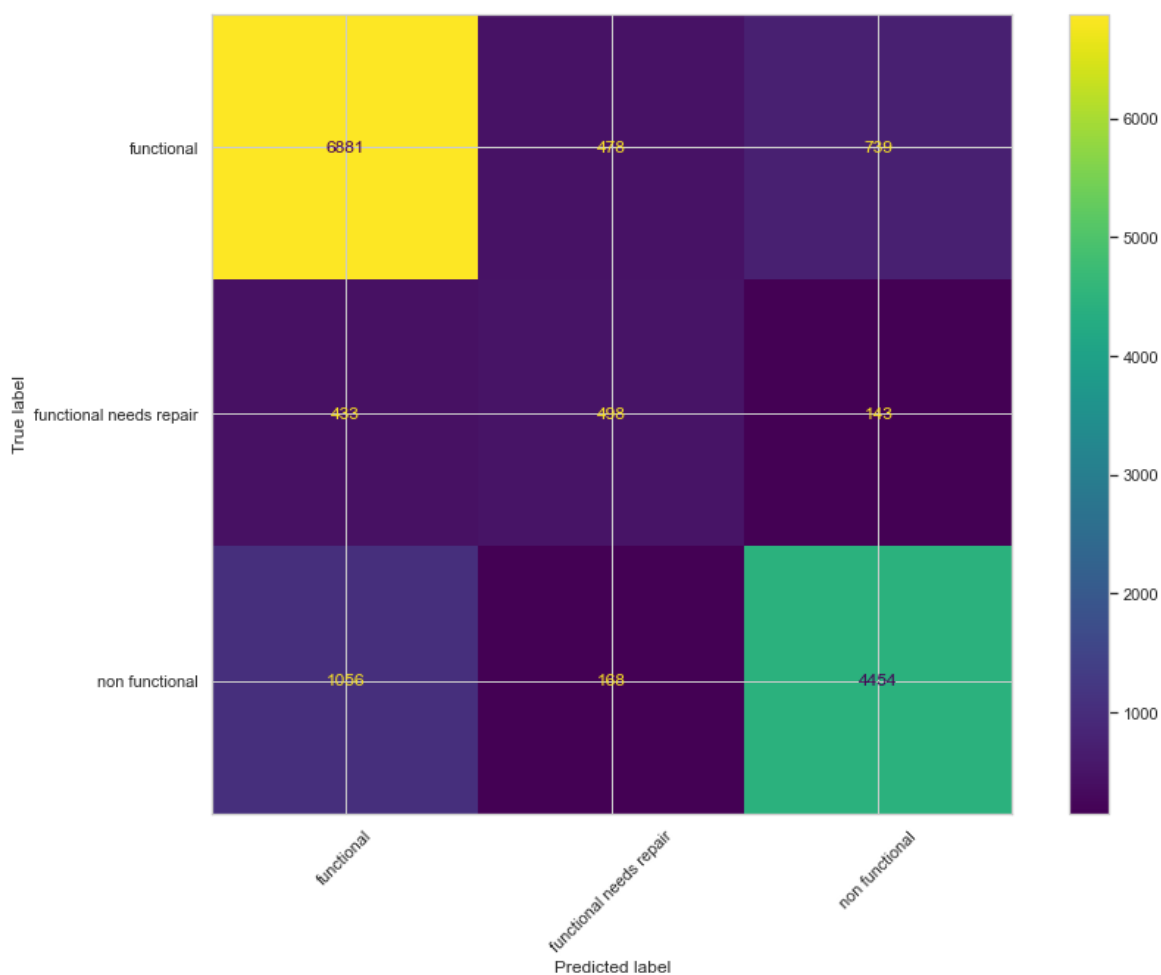
```
In [572]: 1 y_hat_final = final_model.predict(X_test)
```

```
In [573]: 1 print(classification_report(y_test, y_hat_final))
```

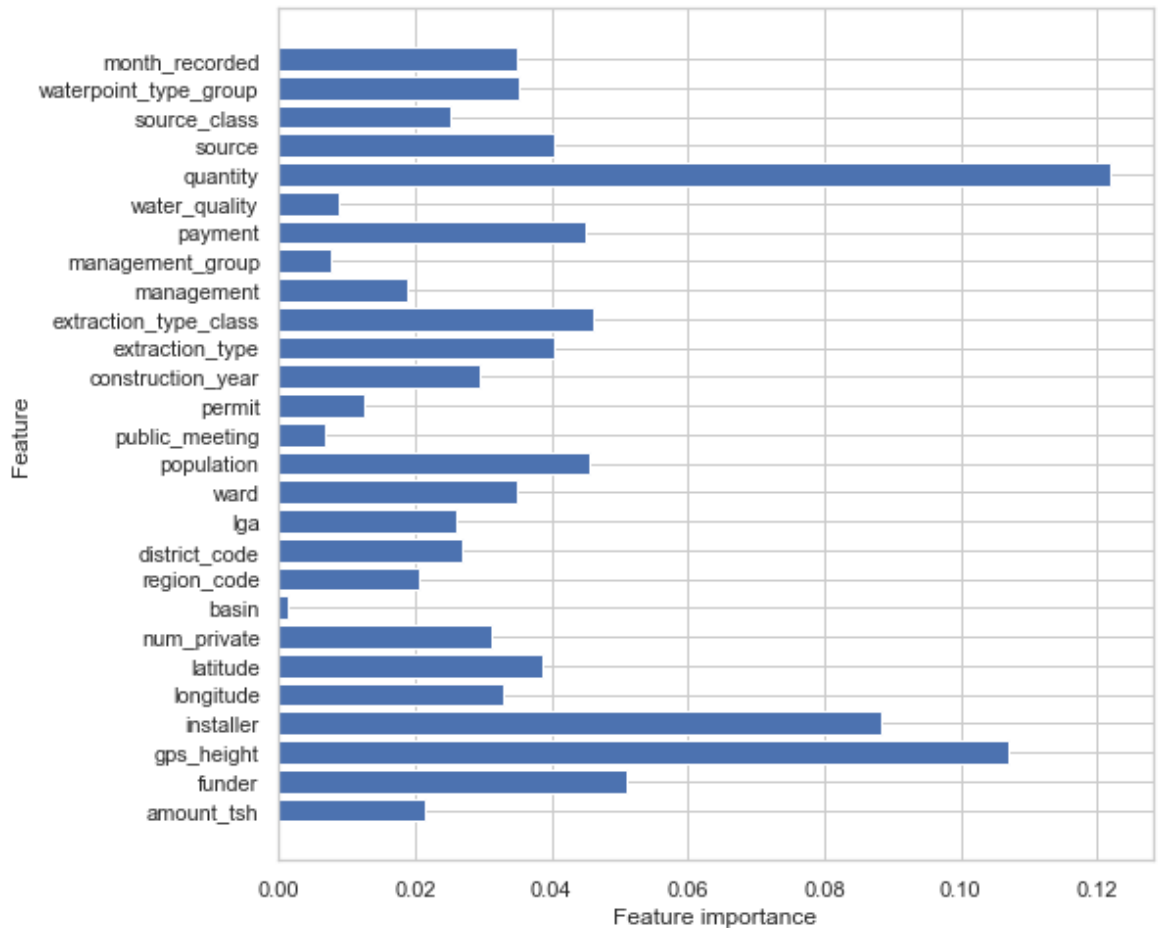
	precision	recall	f1-score	support
functional	0.82	0.85	0.84	8098
functional needs repair	0.44	0.46	0.45	1074
non functional	0.83	0.78	0.81	5678
accuracy			0.80	14850
macro avg	0.70	0.70	0.70	14850
weighted avg	0.80	0.80	0.80	14850

```
In [574]: 1 plot_confusion_matrix(final_model, X_test, y_test,
2                                     xticks_rotation = 45)
```

```
Out[574]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7fc7aecf9bb0>
```



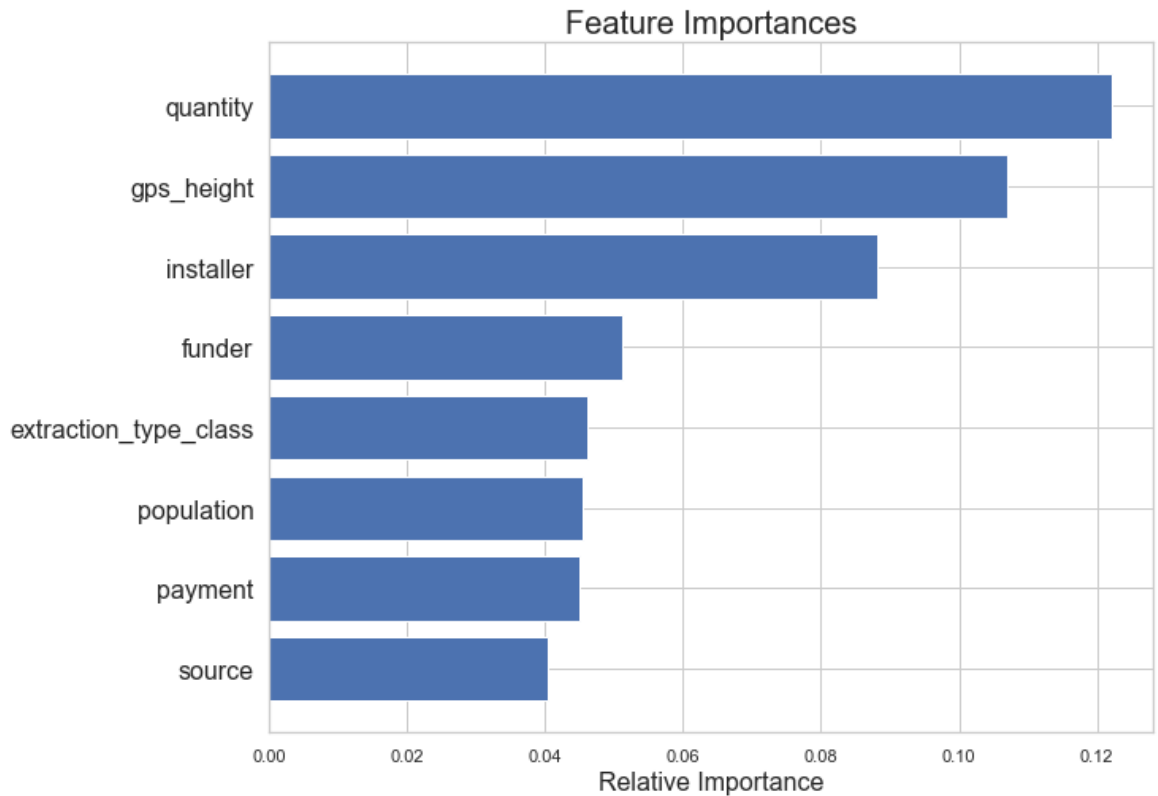
```
In [575]: 1 def plot_feature_importances(model):  
2     n_features = X_train.shape[1]  
3     plt.figure(figsize=(8,8))  
4     plt.barh(range(n_features), model.feature_importances_, align=  
5     plt.yticks(np.arange(n_features), X_train.columns.values)  
6     plt.xlabel('Feature importance')  
7     plt.ylabel('Feature')  
8  
9     plot_feature_importances(final_model.named_steps['rfc'])
```



```

In [620]: 1 # sort the top 8 features
          2
          3 features = X_train.columns
          4 importances = final_model.named_steps['rfc'].feature_importances_
          5 indices = np.argsort(importances)
          6
          7 # customized number
          8 num_features = 8
          9
         10 plt.figure(figsize=(10, 8))
         11 plt.title('Feature Importances', fontsize = 20)
         12
         13 # only plot the customized number of features
         14 plt.barh(range(num_features), importances[indices[-num_features:]])
         15 plt.yticks(range(num_features), [features[i] for i in indices[-num_
         16 plt.xlabel('Relative Importance', fontsize = 16)
         17
         18 plt.savefig('Feature Importance.png')
         19 plt.show();

```



5. Gradient Boost Model

Compare a default parameter gradient boost model against the RFC final model and select final model.

```
In [588]: 1 # Gradient Boost model
          2 gbc_model_pipe = Pipeline([('ct_loo', CT_loo), ('gbc', GradientBoo
          3
          4 gbc_model_pipe.fit(X_train, y_train)
```

```
Out[588]: Pipeline(steps=[('ct_loo',
                           ColumnTransformer(n_jobs=-1, remainder='passthrough',
                           transformers=[('subpipe_num',
                                         Pipeline(steps=[('num_impute',
                                                             SimpleImputer()),
                                                             ('standardScaler',
                                                             StandardScaler())])),
                                         [0, 2, 4, 5, 12]),
                                         ('subpipe_cat',
                                         Pipeline(steps=[('cat_impute',
                                                             SimpleImputer(fill_value='missing',
                                                             strategy='constant'))],
                                         ('categoricalEncoder',
                                         CategoricalEncoder()))])),
                           [1, 3, 6, 7, 8, 9,
                           10, 11, 13,
                           14, 15, 16, 17, 18, 19, 20,
                           21, 22, 23, 24, 25, 26])]),
                           ('gbc', GradientBoostingClassifier(random_state=42))])
```

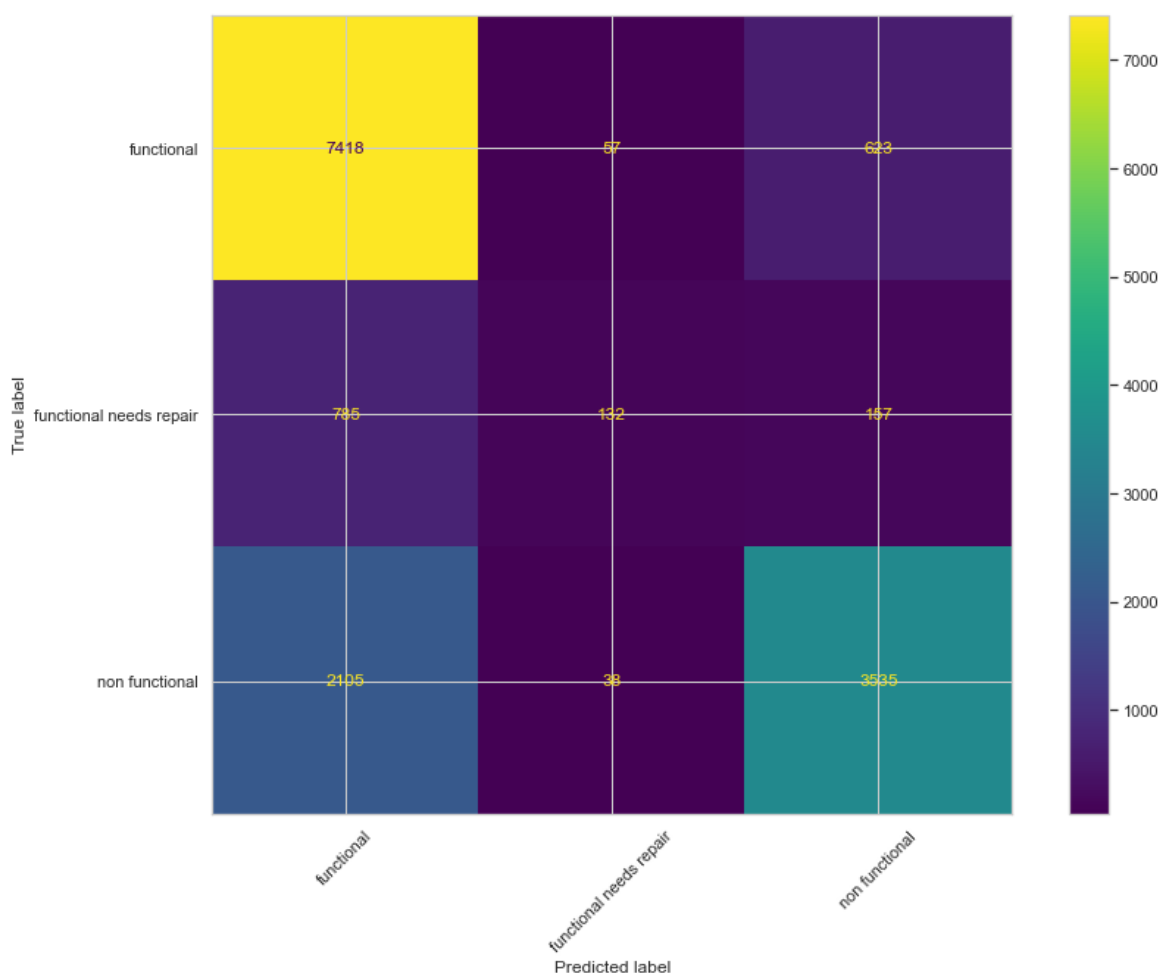
```
In [589]: 1 gbc_model_pipe.score(X_train, y_train)
```

```
Out[589]: 0.7560044893378227
```

```
In [590]: 1 y_hat_gbc = gbc_model_pipe.predict(X_test)
          2
          3 print(classification_report(y_test, y_hat_gbc))
```

	precision	recall	f1-score	support
functional	0.72	0.92	0.81	8098
functional needs repair	0.58	0.12	0.20	1074
non functional	0.82	0.62	0.71	5678
accuracy			0.75	14850
macro avg	0.71	0.55	0.57	14850
weighted avg	0.75	0.75	0.72	14850

```
In [592]: 1 plot_confusion_matrix(gbc_model_pipe, X_test, y_test, xticks_rotat
```



Feature Importance exploration

```
In [621]: 1 # Funder explore
          2 df.funder.value_counts()
```

```
Out[621]: Government Of Tanzania    9084
          Danida                    3114
          Hesawa                    2202
          Rwssp                     1374
          World Bank                1349
          ...
          Mungaya                    1
          Fpct Mulala                1
          Boma Saving                1
          Care Int                   1
          D Ct                       1
          Name: funder, Length: 1897, dtype: int64
```

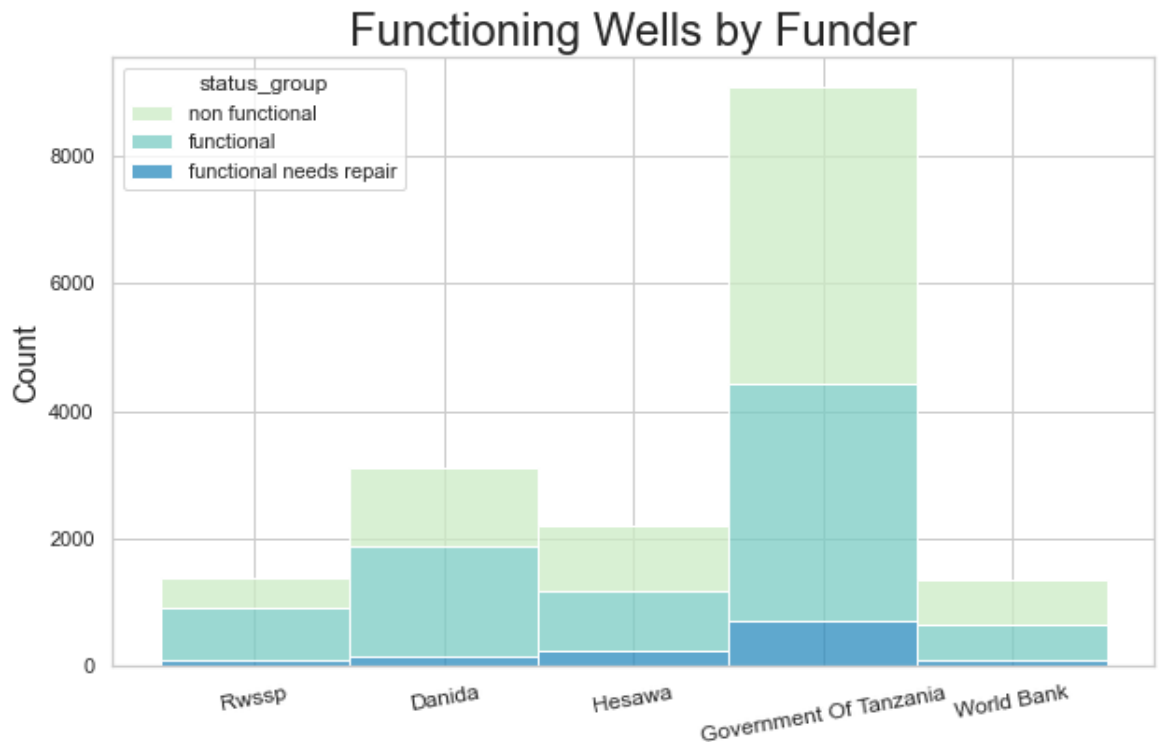
```
In [627]: 1 # Plot four shared cross country water basins in the region
          2 # Plot visual showing basins and functional wells
          3 df_funder = df[df['funder'].isin(['Government Of Tanzania', 'Danida',
          4                                'Rwssp', 'World Bank'
          5                                ])]
          6
```

```
In [683]: 1 df_funder.status_group.value_counts(normalize = True)
```

```
Out[683]: non functional    0.473398
          functional        0.450797
          functional needs repair 0.075804
          Name: status_group, dtype: float64
```



```
In [689]: 1
2 sns.set_theme()
3
4 sns.set(rc={"figure.figsize":(10, 6)})
5 sns.set_style('whitegrid')
6
7 sns.histplot(data = df_funder, x = 'funder', hue = 'status_group',
8             bins = 10, binwidth = 6, palette = 'GnBu', legend = '
9             multiple = 'stack')
10
11
12 plt.title("Functioning Wells by Funder", fontsize= 24)
13 plt.xlabel(None)
14 plt.ylabel("Count", fontsize = 16)
15 plt.xticks(rotation = 10, fontsize = 12)
16
17 plt.savefig('functioning wells by funder.png')
18
19 plt.show();
20
```



```
In [631]: 1 # Explore Installer
          2 df.installer.value_counts()
```

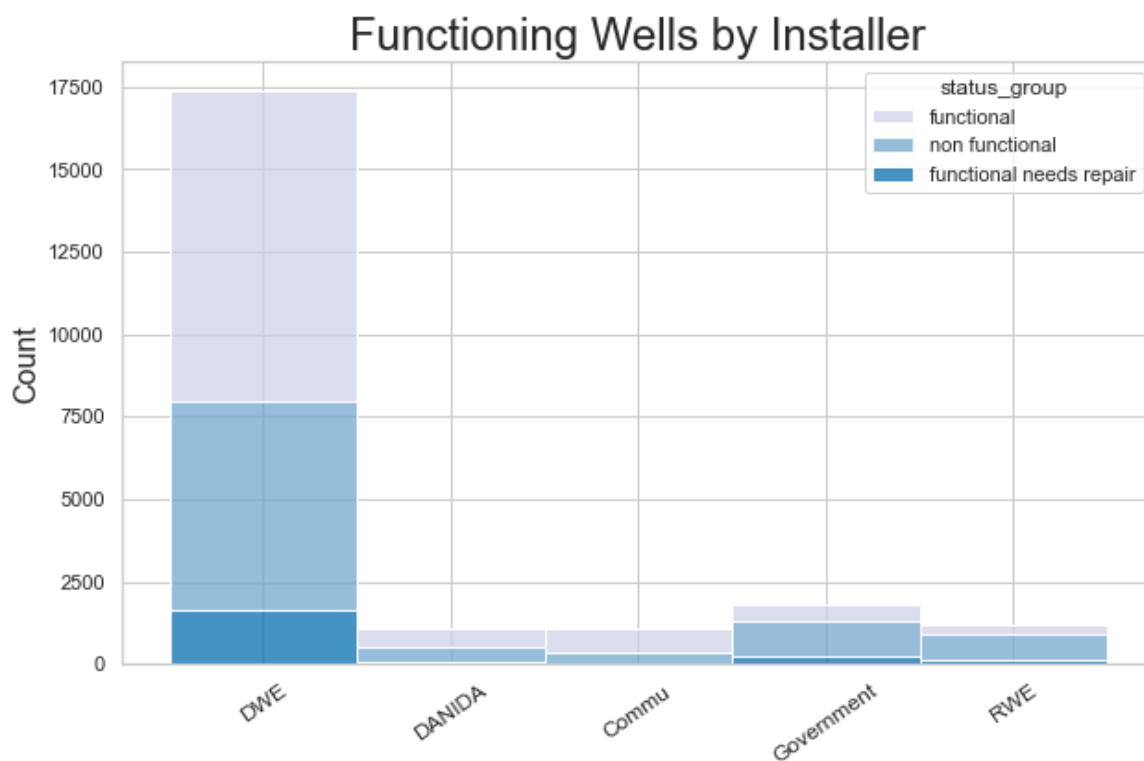
```
Out[631]: DWE                17402
Government                1825
RWE                      1206
Commu                    1060
DANIDA                   1050
...
villager                  1
GACHUMA CONSTRUCTION      1
Atlas                     1
TWESA/JAMII               1
NDM                       1
Name: installer, Length: 2145, dtype: int64
```

```
In [632]: 1 df_installer = df[df['installer'].isin(['Government', 'DANIDA', 'R
          2                                           'Commu', 'DWE'
          3                                           ])]
```

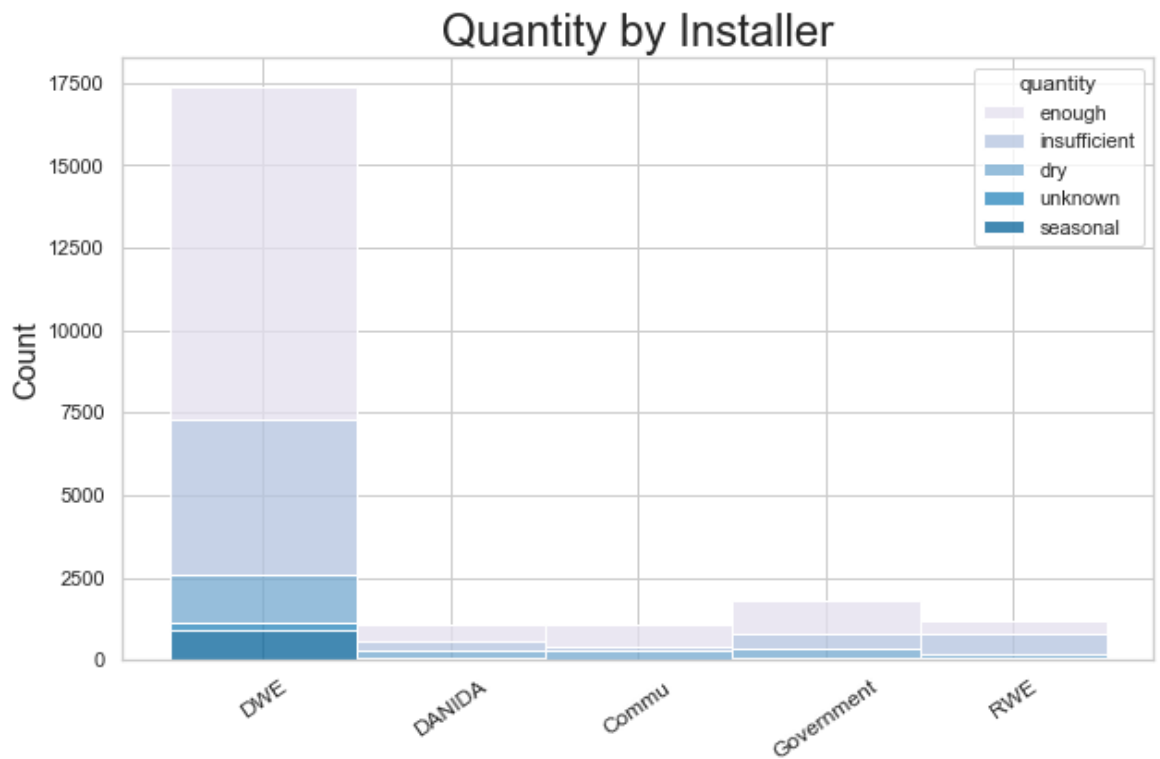
```
In [640]: 1 df_installer.status_group.value_counts(normalize = True)
```

```
Out[640]: functional                0.511822
non functional                      0.393692
functional needs repair             0.094486
Name: status_group, dtype: float64
```

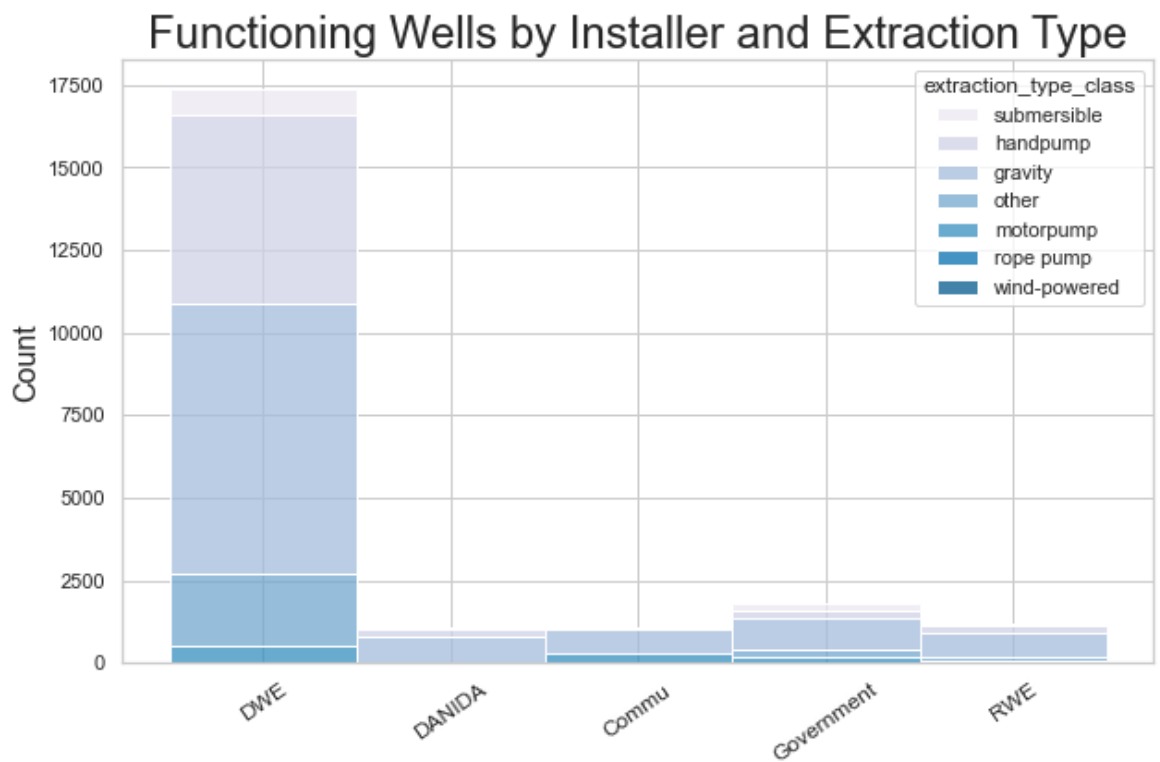
```
In [636]: 1 sns.histplot(data = df_installer, x = 'installer', hue = 'status_g
2           bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3           multiple = 'stack')
4
5
6 plt.title("Functioning Wells by Installer", fontsize= 24)
7 plt.xlabel(None)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 plt.savefig('functioning wells by installer.png')
12
13 plt.show();
14
```



```
In [637]: 1 sns.histplot(data = df_installer, x = 'installer', hue = 'quantity',
2               bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3               multiple = 'stack')
4
5
6 plt.title("Quantity by Installer", fontsize= 24)
7 plt.xlabel(None)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 #plt.savefig('')
12
13 plt.show();
14
```



```
In [638]: 1 sns.histplot(data = df_installer, x = 'installer', hue = 'extracti
2           bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3           multiple = 'stack')
4
5
6 plt.title("Functioning Wells by Installer and Extraction Type", fon
7 plt.xlabel(None)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 #plt.savefig('')
12
13 plt.show();
14 #
15
16
```



```
In [664]: 1 # Sea Level wells
          2 df.gps_height.value_counts(bins = 1000, normalize = True)
```

```
Out[664]: (-1.34, 1.52]          0.345017
          (-15.64, -12.78]      0.002795
          (-18.5, -15.64]       0.002492
          (-21.36, -18.5]       0.002340
          (1288.52, 1291.38]    0.002104
          ...
          (2529.76, 2532.62]    0.000000
          (2526.9, 2529.76]     0.000000
          (2524.04, 2526.9]     0.000000
          (2515.46, 2518.32]    0.000000
          (2572.66, 2575.52]    0.000000
          Name: gps_height, Length: 1000, dtype: float64
```

```
In [641]: 1 # Consider worldbank, top installers, and extraction type
          2 df_worldbank = df[df['funder'].isin(['World Bank'])]
          3
          4 df_worldbank.status_group.value_counts(normalize = True)
          5
```

```
Out[641]: non functional      0.524092
          functional          0.404003
          functional needs repair 0.071905
          Name: status_group, dtype: float64
```

```
In [642]: 1 df_worldbank.installer.value_counts()
```

```
Out[642]: DWE          152
          World        120
          World Bank    95
          Government    57
          WORLD BANK    46
          ...
          SAXON BUILDING CONTRACTORS 1
          EFAM                      1
          JUINE CO                   1
          Wadeco                     1
          Word bank                   1
          Name: installer, Length: 131, dtype: int64
```

```
In [669]: 1 #Clean installer column for world bank
          2 #df.loc[df["gender"] == "male", "gender"] = 1
          3
          4 df_worldbank.loc[df_worldbank['installer'] == 'World', 'installer']
          5 df_worldbank.loc[df_worldbank['installer'] == 'WORLD BANK', 'insta
          6 df_worldbank.loc[df_worldbank['installer'] == 'Word Bank', 'instal
          7 df_worldbank.loc[df_worldbank['installer'] == 'Word bank', 'instal
          8 df_worldbank.loc[df_worldbank['installer'] == 'world bank', 'insta
          9 df_worldbank.loc[df_worldbank['installer'] == 'Word', 'installer']
          10 df_worldbank.loc[df_worldbank['installer'] == 'world', 'installer']
```

```
In [671]: 1 df_worldbank.installer.value_counts(normalize = True)
```

```
Out[671]: World Bank      0.201046
DWE      0.113602
Government 0.042601
Water board 0.023916
Gwasco L   0.020927
...
Conta      0.000747
Building works engineering Ltd 0.000747
GEN        0.000747
Consultant Engineer 0.000747
EFAM       0.000747
Name: installer, Length: 125, dtype: float64
```

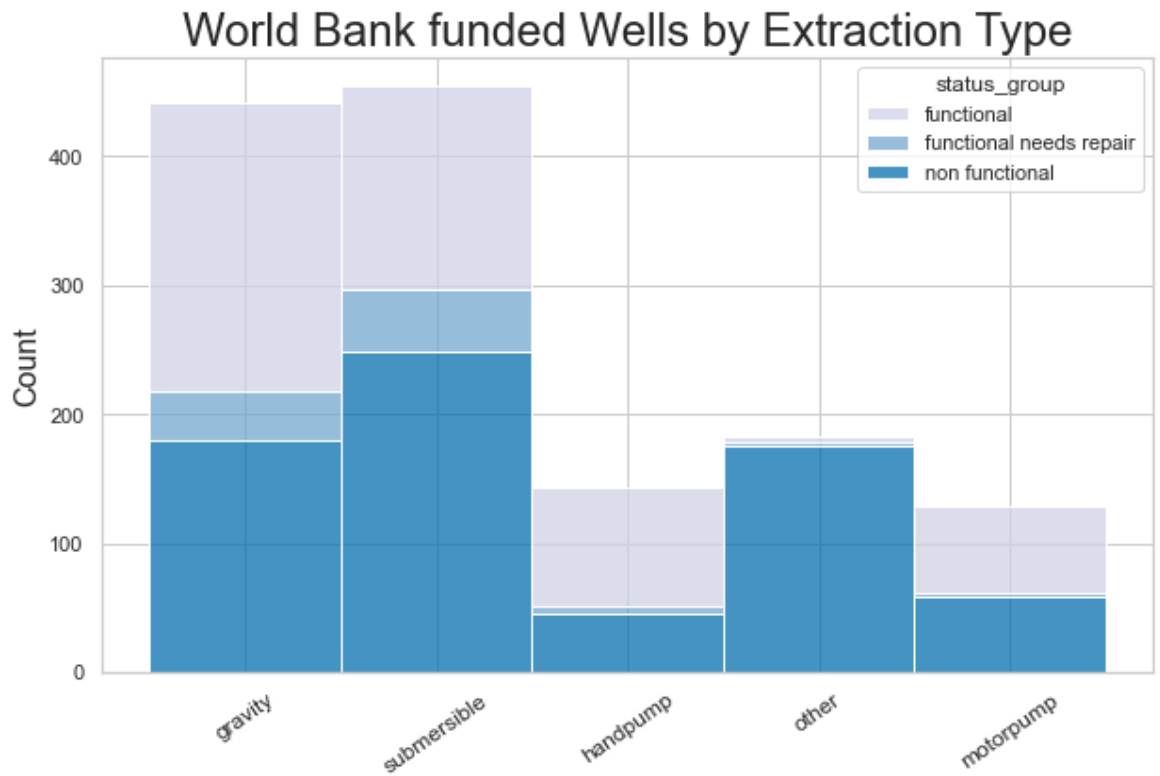
```
In [643]: 1 df_worldbank.extraction_type_class.value_counts(normalize = True)
```

```
Out[643]: submersible    0.336546
gravity    0.326909
other      0.135656
handpump   0.106004
motorpump  0.094885
Name: extraction_type_class, dtype: float64
```

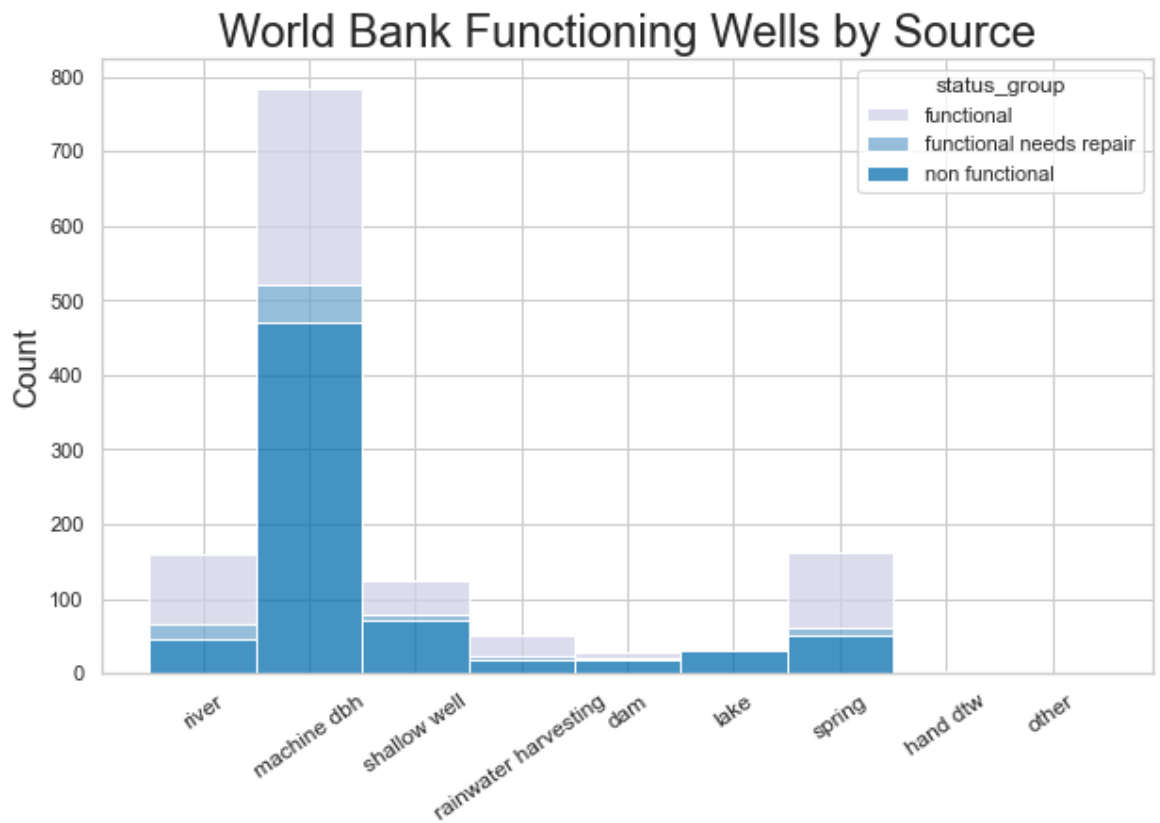
```
In [644]: 1 df_worldbank.source.value_counts(normalize = True)
```

```
Out[644]: machine dbh    0.581913
spring    0.120830
river     0.118606
shallow well 0.092661
rainwater harvesting 0.038547
lake      0.023721
dam       0.020756
hand dtw  0.002224
other     0.000741
Name: source, dtype: float64
```

```
In [666]: 1 sns.histplot(data = df_worldbank, x = 'extraction_type_class', hue
2             bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3             multiple = 'stack')
4
5
6 plt.title("World Bank funded Wells by Extraction Type", fontsize= 2
7 plt.xlabel(None)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 #plt.savefig('')
12
13 plt.show();
14 #
```




```
In [667]: 1 sns.histplot(data = df_worldbank, x = 'source', hue = 'status_group',
2             bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3             multiple = 'stack')
4
5
6 plt.title("World Bank Functioning Wells by Source", fontsize= 24)
7 plt.xlabel(None)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 #plt.savefig('')
12
13 plt.show();
14 #
```



In [672]: 1 df_worldbank.columns

Out[672]: Index(['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'],
dtype='object')

In [674]: 1 *#examine quantity, extraction_type_class, waterpoint_type*
2 df_worldbank.quantity.value_counts()

Out[674]: enough 780
insufficient 239
dry 142
unknown 122
seasonal 66
Name: quantity, dtype: int64

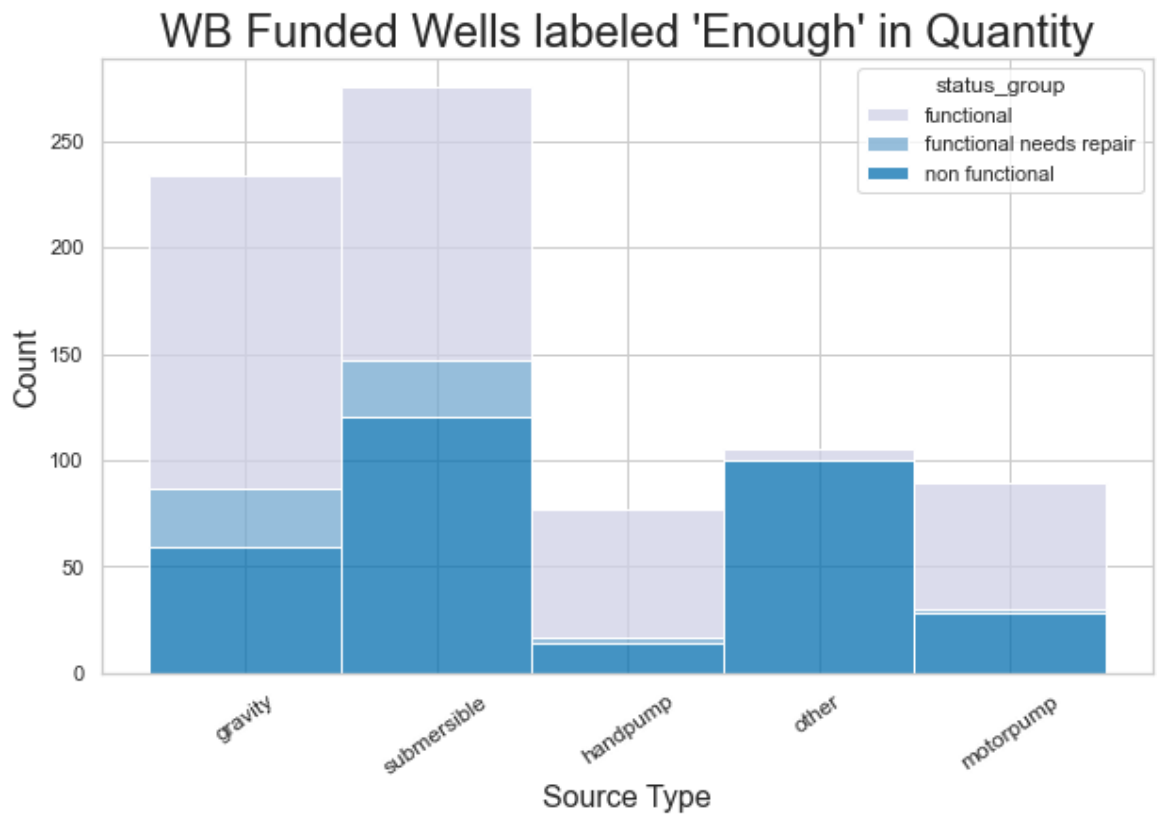
In [687]: 1 *#Examine quantity enough and functioning wells*
2 df_worldbank_enough = df_worldbank[df_worldbank['quantity'].isin(['enough', 'functional'])]
3
4 df_worldbank_enough.status_group.value_counts(normalize = True)

Out[687]: functional 0.511538
non functional 0.411538
functional needs repair 0.076923
Name: status_group, dtype: float64

```

In [686]: 1 sns.histplot(data = df_worldbank_enough, x = 'extraction_type_class',
2                   bins = 10, binwidth = 6, palette = 'PuBu', legend = '
3                   multiple = 'stack')
4
5
6 plt.title("WB Funded Wells labeled 'Enough' in Quantity", fontsize=
7 plt.xlabel("Source Type", fontsize = 16)
8 plt.ylabel("Count", fontsize = 16)
9 plt.xticks(rotation = 35, fontsize = 12)
10
11 plt.savefig('World Bank Enough.png')
12
13 plt.show();
14 #

```



```

In [682]: 1 df_worldbank_enough.extraction_type_class.value_counts(normalize =

```

```

Out[682]: submersible    275
gravity        234
other         105
motorpump      89
handpump       77
Name: extraction_type_class, dtype: int64

```

6. Data Visualizations

Create three data visualizations to communicate findings to Water Aid.

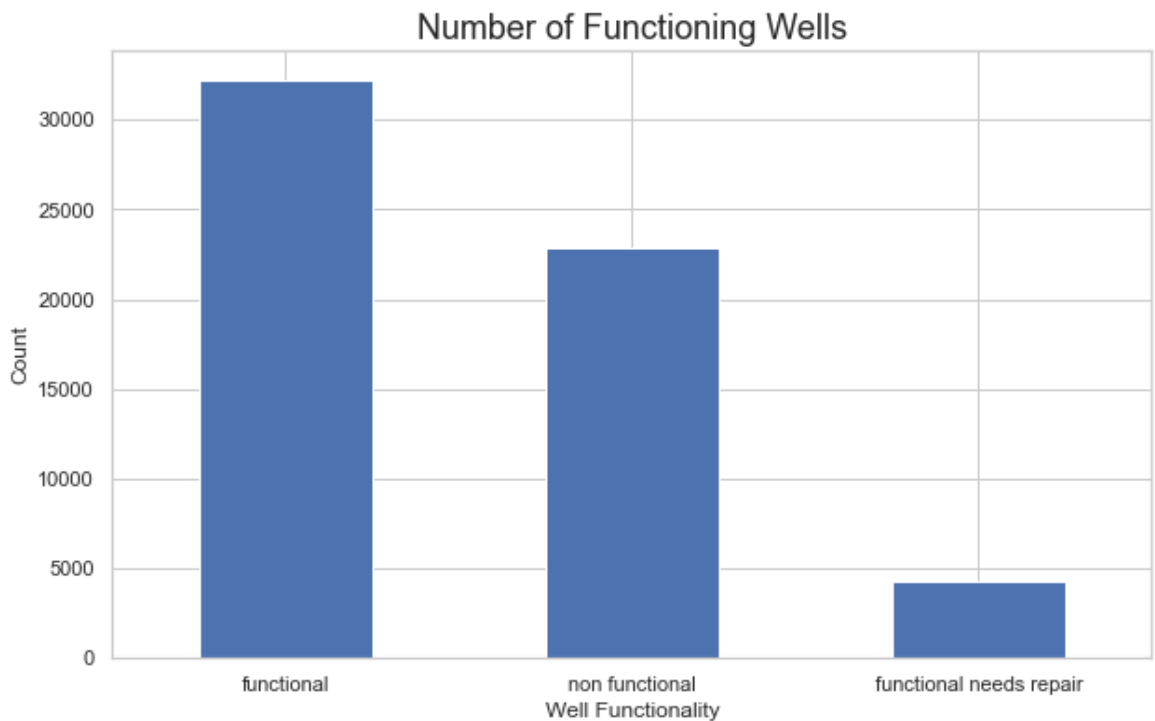
```
In [797]: 1 # Examine full dataframe columns
          2 df.columns
```

```
Out[797]: Index(['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'],
                  dtype='object')
```

```

In [816]: 1 # Visualize well status count from dataset
          2
          3 df.status_group.value_counts().plot(kind="bar")
          4 plt.title("Number of Functioning Wells", fontsize= 18)
          5 plt.xlabel("Well Functionality", fontsize = 12)
          6 plt.xticks(rotation=0)
          7 plt.ylabel("Count", fontsize = 12)
          8
          9 plt.savefig('Number of Functioning Wells.png')
         10
         11 plt.show();
         12
         13

```



```

In [799]: 1 # Plot four shared cross country water basins in the region
          2 # Plot visual showing basins and functional wells
          3 df_basin = df[df['basin'].isin(['Lake Nyasa', 'Lake Victoria', 'La
          4                                         'Ruvuma / Southern Coast',
          5                                         ])]

```

```

In [800]: 1 df_basin.basin.value_counts()

```

```

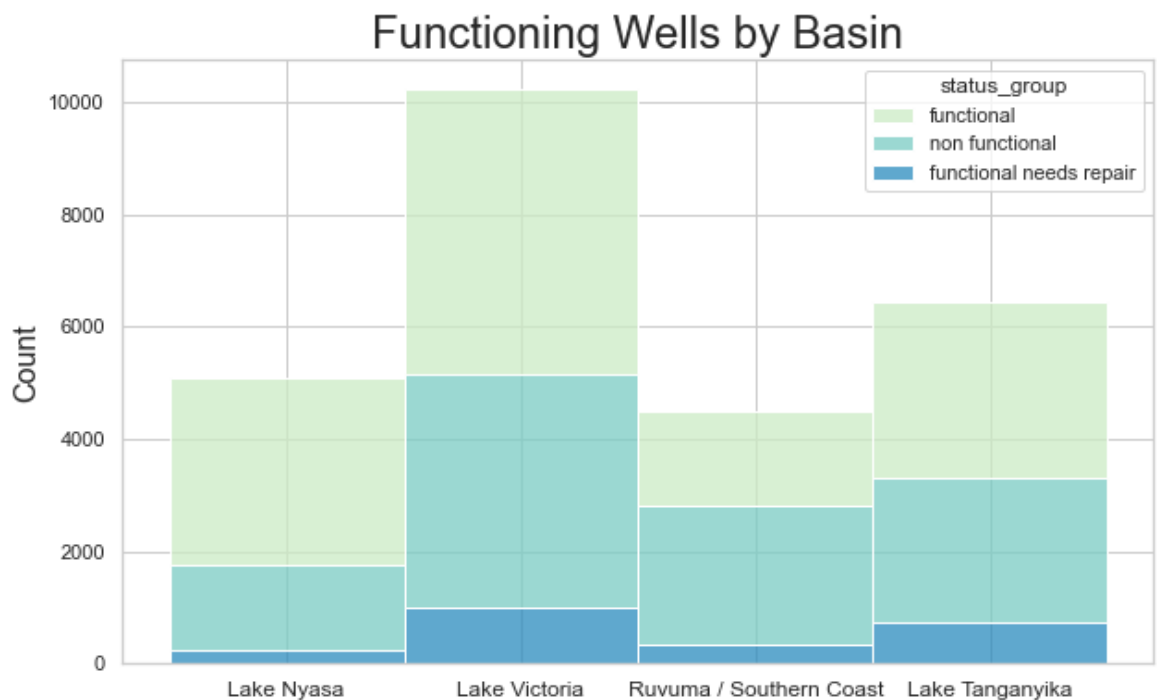
Out[800]: Lake Victoria          10248
          Lake Tanganyika         6432
          Lake Nyasa              5085
          Ruvuma / Southern Coast  4493
          Name: basin, dtype: int64

```

```
In [801]: 1 df_basin.status_group.value_counts()
```

```
Out[801]: functional          13201
non functional          10750
functional needs repair    2307
Name: status_group, dtype: int64
```

```
In [829]: 1 # Use Seaborn to and stacked histogram to show the four basins and
2
3 sns.set_theme()
4
5 sns.set(rc={"figure.figsize":(10, 6)})
6 sns.set_style('whitegrid')
7
8 sns.histplot(data = df_basin, x = 'basin', hue = 'status_group',
9             bins = 10, binwidth = 6, palette = 'GnBu', legend = '
10             multiple = 'stack')
11
12
13 plt.title("Functioning Wells by Basin", fontsize= 24)
14 plt.xlabel(None)
15 plt.ylabel("Count", fontsize = 16)
16 plt.xticks(rotation = 0, fontsize = 12)
17
18 plt.savefig('functioning wells by basin.png')
19
20 plt.show();
21
22
```



```
In [803]: 1 # Examine well status for two basin recommendations: Lake Victoria
          2 df_victoria = df[df['basin'].isin(['Lake Victoria'])]
          3
          4 df_victoria.status_group.value_counts()
```

```
Out[803]: functional          5100
          non functional      4159
          functional needs repair    989
          Name: status_group, dtype: int64
```

```
In [804]: 1 # Create dataframe for Ruvuma Basin
          2 df_ruvuma = df[df['basin'].isin(['Ruvuma / Southern Coast'])]
          3
          4 #Examine well function status for Ruvuma Basin
          5 print(df_ruvuma.status_group.value_counts(normalize = True))
          6 print()
          7 print(df_ruvuma.status_group.value_counts())
```

```
non functional          0.555753
functional              0.371689
functional needs repair  0.072557
Name: status_group, dtype: float64
```

```
non functional          2497
functional              1670
functional needs repair   326
Name: status_group, dtype: int64
```

```
In [805]: 1 df_ruvuma.head()
```

```
Out[805]:
```

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latitude	v
3	67743	0.0	2013-01-28	Unicef	263	UNICEF	38.486161	-11.155298	
26	55012	500.0	2013-01-16	Sobodo	200	Kilolo Star	39.370777	-9.942532	
46	45111	20.0	2013-02-05	Lga	240	LGA	39.087415	-11.000604	M
91	62591	0.0	2013-01-20	Jica	212	Kokeni	38.962945	-10.476566	
98	33379	0.0	2013-02-19	Danida	1000	DWE	35.542173	-10.808853	

5 rows × 41 columns

```
In [806]: 1 #examine descriptive statistics for Ruvuma
          2 df_ruvuma.describe()
```

Out[806]:

	id	amount_tsh	gps_height	longitude	latitude	num_private	region
count	4493.000000	4493.000000	4493.000000	4493.000000	4493.000000	4493.000000	4493.0
mean	37322.257067	228.390385	410.640329	38.316789	-10.485215	0.124861	52.0
std	21489.456338	777.985990	338.566284	1.549237	0.591604	6.745120	38.0
min	19.000000	0.000000	-90.000000	34.889771	-11.649440	0.000000	8.0
25%	18908.000000	0.000000	164.000000	37.244214	-10.850966	0.000000	10.0
50%	37228.000000	0.000000	342.000000	38.935668	-10.626269	0.000000	80.0
75%	55874.000000	50.000000	585.000000	39.448147	-10.250908	0.000000	90.0
max	74247.000000	15000.000000	1641.000000	40.345193	-8.496806	450.000000	99.0

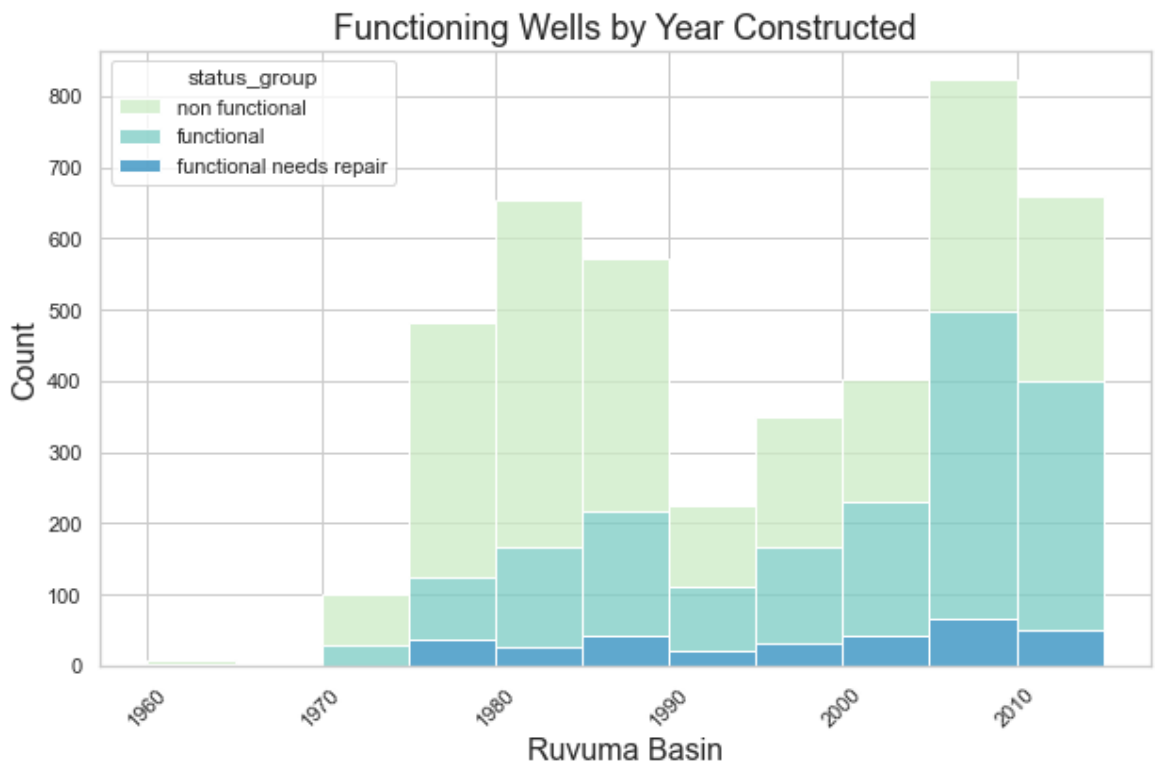
```
In [807]: 1 # Build dataframe ruvuma basin with construction year
          2 df_ruvuma_built = df_ruvuma[df_ruvuma.construction_year != 0]
```



```

In [818]: 1 # Create data visualization, histogram, for functioning wells in R
          2 sns.set_theme()
          3
          4 sns.set(rc={"figure.figsize":(10, 6)})
          5 sns.set_style('whitegrid')
          6
          7 sns.histplot(data = df_ruvuma_built, x = 'construction_year', hue
          8                 bins = 20, binwidth = 5, palette = 'GnBu', legend = '
          9                 multiple = 'stack')
         10
         11
         12 plt.title("Functioning Wells by Year Constructed", fontsize= 18)
         13 plt.xlabel('Ruvuma Basin', fontsize = 16)
         14 plt.xticks(rotation=0)
         15 plt.ylabel("Count", fontsize = 16)
         16 plt.xticks(rotation = 45)
         17
         18 plt.savefig('functioning wells ruvuma.png')
         19
         20 plt.show()
         21
         22 ;

```



Out [818]: ''

```

In [809]: 1 # Dataframe for older wells in Ruvuma, 1975 to 1990
          2 df_ruvuma_old_wells = df_ruvuma.loc[(df_ruvuma.construction_year >
          3                 (df_ruvuma.construction_year <=

```

```
In [810]: 1 # Percentage of wells in need of repair
          2 df_ruvuma_old_wells.status_group.value_counts(normalize = True)

Out[810]: non functional    0.699024
          functional        0.236646
          functional needs repair 0.064331
          Name: status_group, dtype: float64
```

Final Summary

The initial Logistic Regression model with default parameters delivered the following scores:

		precision	recall	f1-score	sup
port					
	functional	0.78	0.88	0.83	80
98					
	functional needs repair	0.55	0.25	0.34	10
74					
	non functional	0.80	0.74	0.77	56
78					
	accuracy			0.78	148
50					

The final Random Forest Classifier model with class imbalance adjustments and hyperparameter tuning delivered the following scores:

		precision	recall	f1-score	sup
port					
	functional	0.82	0.85	0.84	80
98					
	functional needs repair	0.44	0.46	0.45	10
74					
	non functional	0.83	0.78	0.81	56
78					
	accuracy			0.80	148
50					

The Random Forest Classifier was trained using both Randomized Search and Grid Search. Here are the final hyperparameter adjustments:

n_estimators = 100, (default)

max_depth = 100,

max_features = 'auto',

```
min_samples_leaf = 2,
```

```
min_samples_split = 3
```

The f1 scores increased by 1% for functional wells and increased by 4% for non-functional wells. The f1 scores for functioning wells in need of repair increased the most, by 11%, but that score remains too low to be a reliable predictor of well function. One challenge appears to be many wells in need of repair often end up being classified as functioning wells, their features appear to look much like functioning wells.

Water Aid's use of the model will be primarily for precision - true positive identification for functioning and non-functioning wells. They can still make use of the model for wells in need of repair, but the results for that class need to be understood in terms of recall - the ability of a model to find all the relevant cases within a data set. The f1 score may be the most useful in terms of communicating the model's effectiveness.

Recommendations

1. Model Use

Water aid should focus on identifying non-functioning wells rather than wells that are functioning and in need of repair. Chances are 81% that they will be right which will help in making use of programming resources to repair the wells.

2. Funders

Water Aid should work with the World Bank to repair non-functioning wells. The feature importances from the model reveal that funders and installers are important drivers in whether a well is functioning or not. The World Bank is one of the top 5 funders, with whom Water Aid already has a working relationship. World Bank also installs 20% of the wells they fund. Water Aid should work more closely with World Bank, especially considering that 52% of the wells that the World Bank funds are nonfunctional.

3.

Water Aid should focus on wells that are driven by submersible pumps and have 'enough' water. Out of wells funded by World Bank that have been labeled with 'enough' water in quantity, 41% are not functional. Many of these wells are submersible pumps that draw from deep water wells, these pump types require greater investment and maintenance. Many of the other wells are gravity fed pumps, which are often simple in design but vary greatly in ability to deliver water due to siting issues and basic infrastructure resources.

