

```
In [57]: 1 # Import standard packages
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 %matplotlib inline
7 sns.set_theme(style="darkgrid")
8 from collections import Counter
9 from sklearn.preprocessing import OneHotEncoder, LabelEncoder
10 import time
11 # from sklearn import metrics
12 from sklearn.metrics import precision_score, recall_score, accuracy_score
13 from sklearn.pipeline import Pipeline
14 from sklearn.preprocessing import MinMaxScaler, StandardScaler
15 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
16 from sklearn.metrics import confusion_matrix, plot_confusion_matrix
17
18 from sklearn.linear_model import LogisticRegression
19 from sklearn.ensemble import RandomForestClassifier
20 from xgboost import XGBClassifier
```

Tanzania Machine Learning Water Pump Classification

Modeling Notebook

Author: Dylan Dey

This project is available on github here: (insert link here)

The Author can be reached at the following email: ddey2985@gmail.com
(<mailto:ddey2985@gmail.com>)

Classification Metric Understanding

Below is a confusion matrix that would be produced from a model performing predictive maintenance on behalf of the Ministry of Water. There are four possible outcomes to be considered. The confusion matrix below is a visual aid to help in understanding what classification metrics to consider when building the model.

ACTUAL	Predicted	
	Functional	Non-Functional
Functional	True Positive A true positive for my model will be considered a functional pump correctly labeled as functional. This means no teams and/or additional resources will need to be sent here.	False Negatives A false negative for my model would happen when the model predicts a well that is actually functional to be non-functional. Therefore, teams/resources would be sent to communities that already have a functional water pump. Reducing false positives would help efficiency of distributing resources appropriately.
Non-Functional	False Positive A false positive in my model would be a non-functional pump incorrectly labeled as a functional pipe. False positives should be reduced as much as possible as this is the worst case scenario for my model. Teams/resources would be withheld from communities that need them.	True Negatives A true negative for my model would be a non-functional pump correctly labeled as a non-functional pump. Teams and/or additional resources will need to be sent here first.

A true positive in the current context would be when the model correctly identifies a functional pump as functional. A true negative would be when the model correctly identifies a non-functional pump as non-functional. Both are important and both can be described by the overall **accuracy** of the model.

True negatives are really at the heart of the model, as this is the situation in which the Ministry of Water would have a call to action. An appropriately outfitted team would be set to *all* pumps that my model identifies as non-functional. Thus, this is the situation in which the correct resources are being derived to the correct water pumps as quickly as possible. High accuracy would mean that more resources are going to the correct locations from the get-go.

True positives are also important. This is where the model will really be saving time, resources, and money for the Ministry of Water. *Any* pumps identified as functional would no longer need to be physically checked and the Ministry of Water can withhold additional resources from going to pumps that do not actually need them.

Notice the emphasis on *any* and *all* pumps in my description of true negatives and true positives above. The true cost/resource analysis is really the consideration of this fact: no model I create will ever correctly identify every single pump appropriately. This is the cost of predictive maintenance and a proper understanding of false positives and false negatives is extremely important in production of classification models in the given context.

False positives in the current context are the worst case scenario for modeling. This is the scenario in which the model **incorrectly** identifies a non-functional model as functional. Thus, resources would be withheld and no team would be sent to physically check these pumps, as the Ministry of Water would have to assume they are indeed functional if they want to use the model appropriately. False positives therefore describe the number of non-functional pumps that will go unvisited and unfixed until they can be resolved by other means. Reducing false positives as much as possible is very important.

Well, why would I want to build the model if these false positives cannot be completely avoided? Cost/resource management, of course! Afterall, it is about making sure as many people get clean water as quickly as possible. The reality is there that resources are finite, and without the model the Ministry of Water likely would not have the resources to physically check all the pumps and then fix all of the pumps in any sort of reasonable timeline, and even less communities would have access to fresh water when compared to using the model for predictive maintenance.

False negatives are also important to consider. While false positives can be considered more harmful overall, false negatives are also important to reduce as much as possible. In the given context, false negatives describe the situation in which the model **incorrectly** identifies a functional pump as non-functional. Because the Ministry of water will deploy fully equipped teams to visit all pumps that my model predicts to be non-functional, these will be the pumps that will waste resources. Resources will be sent to locations that they aren't needed, and the metric that describes this would be false negatives. Thus, reduction of false negatives is essential in improving the efficiency of resource management through predictive maintenance.

In summary, overall accuracy of the model and a reduction of both false negatives and false positives are the most important metrics to consider when developing a model in this context. More specifically, models will be tuned to **maximize accuracy and f1-score**.

Accuracy:

f1-score:

Function Definition

Below are all of the functions used for preprocessing data before modeling.

```

In [2]: 1 def drop_cols(water_pump_df):
2         to_drop_final = ['id', 'recorded_by', 'num_private',
3                           'waterpoint_type_group', 'source',
4                           'source_class', 'extraction_type',
5                           'extraction_type_group', 'payment_type',
6                           'management_group', 'scheme_name',
7                           'water_quality', 'quantity_group',
8                           'scheme_management', 'longitude',
9                           'latitude', 'date_recorded',
10                          'amount_tsh', 'gps_height',
11                          'region_code', 'district_code']
12         #'population'
13
14         return water_pump_df.drop(columns=to_drop_final, axis=1)
15
16         #helper function to bin construction year
17 def construction_wrangler(row):
18     if row['construction_year'] >= 1960 and row['construction_year'] <
19         return '60s'
20     elif row['construction_year'] >= 1970 and row['construction_year']
21         return '70s'
22     elif row['construction_year'] >= 1980 and row['construction_year']
23         return '80s'
24     elif row['construction_year'] >= 1990 and row['construction_year']
25         return '90s'
26     elif row['construction_year'] >= 2000 and row['construction_year']
27         return '00s'
28     elif row['construction_year'] >= 2010:
29         return '10s'
30     else:
31         return 'unknown'
32
33 def bin_construction_year(water_pump_df):
34     water_pump_df['construction_year'] = water_pump_df.apply(lambda row:
35     return water_pump_df
36
37
38     #takes zero placeholders and NAN values and converts them into 'unknown'
39 def fill_unknowns(water_pump_df):
40     installer_index_0 = water_pump_df['installer'] == '0'
41     funder_index_0 = water_pump_df['funder'] == '0'
42     water_pump_df.loc[installer_index_0, 'installer'] = 'unknown'
43     water_pump_df.loc[funder_index_0, 'funder'] = 'unknown'
44     water_pump_df.fillna({'installer': 'unknown',
45                           'funder': 'unknown',
46                           'subvillage': 'unknown'}, inplace=True)
47     return water_pump_df
48
49     #returns back boolean features without NANs while maintaining same rat
50 def fill_col_normal_data(water_pump_df):
51     filt = water_pump_df['permit'].isna()
52     probs = water_pump_df['permit'].value_counts(normalize=True)
53     water_pump_df.loc[filt, 'permit'] = np.random.choice([True, False]
54                                                           size=int(filt.sum()),
55                                                           p = [probs[True], probs[False]])
56     filt = water_pump_df['public_meeting'].isna()

```

```

57     probs = water_pump_df['public_meeting'].value_counts(normalize=True)
58     water_pump_df.loc[filt, 'public_meeting'] = np.random.choice([True
59         size=int(filt.sum()),
60         p = [probs[True], probs[False]])
61     return water_pump_df
62
63
64
65 def apply_cardinality_reduct(water_pump_df, reduct_dict):
66     for col, categories_list in reduct_dict.items():
67         water_pump_df[col] = water_pump_df[col].apply(lambda x: x if x
68     return water_pump_df
69
70
71
72 #one_hot_incode categorical data
73 def one_hot(water_pump_df):
74     final_cat = ['funder', 'installer', 'wpt_name', 'basin', 'subvilla
75         'lga', 'ward', 'public_meeting', 'permit', 'construction_year',
76         'extraction_type_class', 'management', 'payment', 'quality_grou
77         'quantity', 'source_type', 'waterpoint_type']
78
79     water_pump_df = pd.get_dummies(water_pump_df[final_cat], drop_firs
80
81     return water_pump_df
82
83
84
85 #master function for cleaning dataframe
86 def clean_dataframe(water_pump_df, reduct_dict):
87     water_pump_df = drop_cols(water_pump_df)
88     water_pump_df = bin_construction_year(water_pump_df)
89     water_pump_df = fill_unknowns(water_pump_df)
90     water_pump_df = fill_col_normal_data(water_pump_df)
91     water_pump_df = apply_cardinality_reduct(water_pump_df, reduct_dic
92     water_pump_df = one_hot(water_pump_df)
93
94     return water_pump_df
95
96 #####
97 # The rest of the functions in this section
98 #define functions that reduce cardinality
99 #by mapping infrequent values ot other
100 #the dictionary derived from these functions
101 #will be used by my_funk in my master
102 #clean_dataframe function
103
104 #helper function for reducing cardinality
105 def cardinality_threshold(column, threshold=0.65):
106     #calculate the threshold value using
107     #the frequency of instances in column
108     threshold_value=int(threshold*len(column))
109     #initialize a new list for lower cardinality column
110     categories_list=[]
111     #initialize a variable to calculate sum of frequencies
112     s=0
113     #Create a dictionary (unique_category: frequency)

```

```

114     counts=Counter(column)
115
116     #Iterate through category names and corresponding frequencies after
117     #by descending order of frequency
118     for i,j in counts.most_common():
119         #Add the frequency to the total sum
120         s += dict(counts)[i]
121         #append the category name to the categories list
122         categories_list.append(i)
123         #Check if the global sum has reached the threshold value, if s
124         if s >= threshold_value:
125             break
126         #append the new 'Other' category to list
127         categories_list.append('Other')
128
129     #Take all instances not in categories below threshold
130     #that were kept and lump them into the
131     #new 'Other' category.
132     new_column = column.apply(lambda x: x if x in categories_list else
133     #     return new_column
134     return categories_list
135
136     #reduces the cardinality of appropriate categories
137 def get_col_val_mapping(water_pump_df):
138     col_threshold_list = [
139         ('funder', 0.65),
140         ('installer', 0.65),
141         ('wpt_name', 0.15),
142         ('subvillage', 0.07),
143         ('lga', 0.6),
144         ('ward', 0.05)
145     ]
146
147     reduct_dict = {}
148
149     for col, thresh in col_threshold_list:
150         reduct_dict[col] = cardinality_threshold(water_pump_df[col],
151         threshold= thresh)
152
153     return reduct_dict
154
155     # reduct_dict is a key value mapper that will
156     # be used for both training and testing sets
157     # in order to reduce cardinality of the data

```

Import The Data From Multiple Sources

I used a number of sources for my data to use for modeling in this notebook. The cell below imports the original data from the DrivenData competition, data derived from DrivenData in QGIS and opensource hydrology data, and population data from Tanzania government census in 2012.

```

In [3]: 1 #import data from DrivenData
2 train_labels = pd.read_csv('files/0bf8bc6e-30d0-4c50-956a-603fc693d966.
3 train_features = pd.read_csv('files/4910797b-ee55-40a7-8668-10efd5c1b96
4 df = train_features.merge(train_labels, on='id').copy()
5 #import QGIS derived data and prepare for model
6 river_df = pd.read_csv('data/river_dist2.csv')
7 #removing outliers
8 index_riv = river_df[river_df['HubDist'] > 66].index
9 river_median = river_df['HubDist'].median()
10 river_df.loc[index_riv, 'HubDist'] = river_median
11
12 #create boolean for pump being within 8 km of river
13 river_s = river_df['HubDist'].copy()
14 river_s.rename('near_river', inplace=True)
15 near_river = river_s[river_s < 8].apply(lambda x: 1 if not pd.isnull(x)
16 df = df.join(near_river)
17 df.near_river.fillna(0, inplace=True)
18
19 #import population data from 2012 government census
20 df_pop = pd.read_excel('data/tza-pop-popn-nbs-baselinedata-xlsx-1.xlsx')
21
22 #create a dictionary of values with format {Ward : Total Population}
23
24 pop_index = df_pop.groupby('Ward_Name')['total_both'].sum().index
25 pop_values = df_pop.groupby('Ward_Name')['total_both'].sum().values
26 pop_dict = dict(zip(pop_index, pop_values))
27
28 #create pandas Dataframe for merging
29 pop_dataframe = pd.DataFrame.from_dict(pop_dict, orient='index')
30 #rename column for clarity
31 pop_dataframe.rename(columns={0: 'ward_pop'}, inplace=True)
32
33 #merge dataframes
34 df_pop_merge = df.merge(pop_dataframe,
35                          how='left',
36                          left_on='ward',
37                          right_index=True)
38
39 #replace null values of ward population with
40 #median ward population
41
42 ward_pop_median = df_pop_merge['ward_pop'].median()
43 df_pop_merge.fillna(value=ward_pop_median, inplace=True)
44 # merge back into df and drop pop column
45 ward_pop_s = df_pop_merge['ward_pop'].copy()
46 df = df.join(ward_pop_s)
47 df.drop(columns=['population'], axis=1, inplace=True)
48
49
50 need_repair_index = df['status_group'] == 'functional needs repair'
51 df_binary = df.copy()
52 df_binary.loc[need_repair_index, 'status_group'] = 'non functional'

```

Transform Target to Numerical Data with Label Encoding

As shown at the beginning of the project, the functional state of the water supply point is described as: functional — It is working; non functional — it is not working; functional needs repair — Running, but needing maintenance.

These groups will be relabeled as: functional: 0 functional needs repair: 1 non functional: 2

This would be for future work in which I didn't bin functional needs repair with non functional. I chose this simplified way to deal with the imbalanced dataset not only due to time-constraints but because it is extremely difficult to get meaningful predictive power for the 'functional needs repair' label due to the nature of the dataset.

For the modeling in this notebook the target data will be relabeled as the following after binning function needs repair with non-functional:

functional: 0

non_functional: 1

```
In [4]: 1 df_binary['status_group'].value_counts(normalize=True)
```

```
Out[4]: functional      0.543081
non functional    0.456919
Name: status_group, dtype: float64
```

```
In [5]: 1 # le = LabelEncoder()
2
3 # le.fit(['functional', 'functional needs repair', 'non functional'])
4
5 # df['status_group'] = le.transform(df.status_group)
6
7 # df['status_group'].value_counts(normalize=True)
```

```
In [6]: 1 le = LabelEncoder()
2
3 le.fit(['non functional', 'functional'])
4
5 df_binary['status_group'] = le.transform(df_binary.status_group)
6
7 df_binary['status_group'].value_counts(normalize=True)
```

```
Out[6]: 0    0.543081
1    0.456919
Name: status_group, dtype: float64
```

Split Data into Test and Training Sets

```
In [7]: 1 from sklearn.model_selection import train_test_split
2
3 y = df_binary['status_group']
4 X = df_binary.drop(columns=['status_group'], axis=1)
5
6 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.3,
```


Run Master Cleaning Function

```
In [8]: 1 value_maps = get_col_val_mapping(df_binary)
        2 X_train = clean_dataframe(X_train, value_maps)
        3 X_test = clean_dataframe(X_test, value_maps)
        4 display(X_train.shape)
        5 X_test.shape
```

(41580, 239)

Out[8]: (17820, 239)

Data Modeling

Baseline model

I decided to use logistic regression as my baseline model. I chose to use an arbitrary large value for C and set the solver to 'liblinear.' I did not choose to fit an intercept for the baseline model.

```
In [9]: 1 # Instantiate the model
2 logreg = LogisticRegression(fit_intercept=False, C=1e12, solver='liblin
3
4 # Fit the model
5 logreg.fit(X_train, y_train)
6 y_hat_train = logreg.predict(X_train)
7 y_hat_test = logreg.predict(X_test)
8
9 print('Training Precision: ', precision_score(y_train, y_hat_train))
10 print('Testing Precision: ', precision_score(y_test, y_hat_test))
11 print('\n\n')
12
13 print('Training Recall: ', recall_score(y_train, y_hat_train))
14 print('Testing Recall: ', recall_score(y_test, y_hat_test))
15 print('\n\n')
16
17 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train))
18 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test))
19 print('\n\n')
20
21 print('Training F1-Score: ', f1_score(y_train, y_hat_train))
22 print('Testing F1-Score: ', f1_score(y_test, y_hat_test))
```

Training Precision: 0.7951634068612136
Testing Precision: 0.7886786918903065

Training Recall: 0.6681543712260436
Testing Recall: 0.6642786561264822

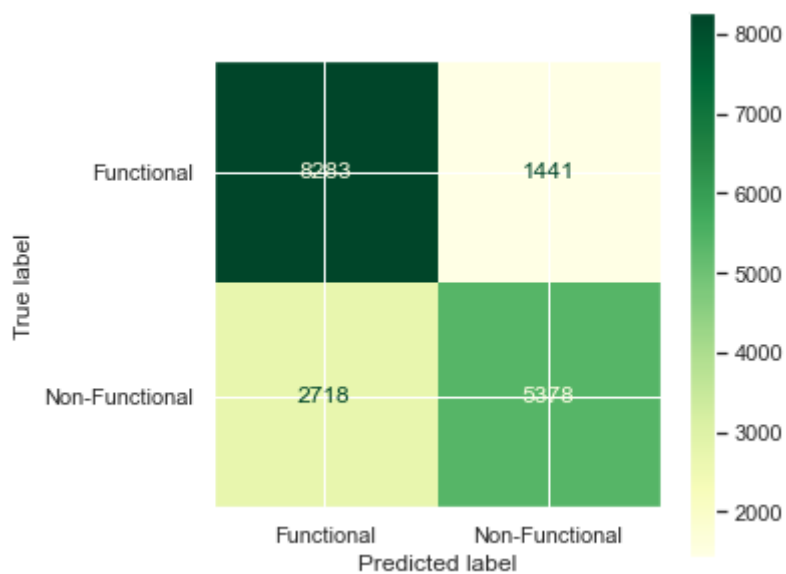
Training Accuracy: 0.7691678691678692
Testing Accuracy: 0.7666105499438832

Training F1-Score: 0.7261469984021912
Testing F1-Score: 0.7211532014750252

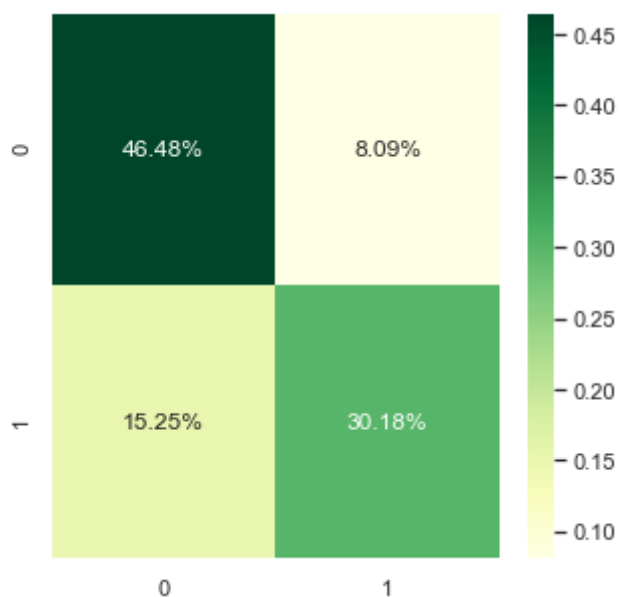
```
In [10]: 1 #Get the confusion matrix
2 lr_base_matrix = confusion_matrix(y_test, y_hat_test)
3 print(lr_base_matrix)
```

```
[[8283 1441]
 [2718 5378]]
```

```
In [71]: 1 fig_labels = ['Functional', 'Non-Functional']
2 fig, ax = plt.subplots(figsize=(5, 5))
3
4 plot_confusion_matrix(logreg,
5                       X_test,
6                       y_test,
7                       ax=ax,
8                       cmap='YlGn',
9                       display_labels=fig_labels)
10
11 plt.show();
```



```
In [70]: 1 # Visualize your confusion matrix
2 fig, ax = plt.subplots(figsize=(5, 5))
3
4
5 sns.heatmap(lr_base_matrix/np.sum(lr_base_matrix), annot=True,
6             fmt='.2%', cmap='YlGn', ax=ax)
7
8 plt.show();
```



Below a random grid search was used to create 300 of 600 possible combinations to narrow down the parameters to do a full brute force gridCV search. The score used for determining the 'best' model was precision.

This is clearly not the way to build the model. Although I am greatly concerned with increasing my precision in order to reduce false positives, false negatives still have a fairly high cost as well.

The precision of the model is excellent. However, the accuracy and F1 scores are unacceptable.

I will describe why this is a poor model in more detail below under the display of the confusion matrix.

```
In [20]: 1 start = time.time()
2
3 classifier_penalties = ['l1', 'l2']
4 classifier_Cs = np.logspace(-5, 5, 100)
5 classifier__solver = ['liblinear']
6
7 random_grid = {'penalty' : classifier_penalties,
8 'C' : classifier_Cs,
9 'solver' : classifier__solver}
10
11 lr = LogisticRegression()
12 # Random search of parameters, using 3 fold cross validation,
13 # search across 100 different combinations, and use all available cores
14 lr_random_precision = RandomizedSearchCV(scoring='precision',
15 estimator = lr,
16 param_distributions = random_grid,
17 n_iter = 100, cv = 3, verbose=2,
18 random_state=42,
19 n_jobs = -1)
20 # Fit the random search model
21 lr_random_precision.fit(X_train, y_train)
22
23
24
25 end = time.time()
26 print(end - start)
27
28 lr_best_rand_precision = lr_random_precision.best_estimator_
29 lr_best_rand_precision.fit(X_train,y_train)
30 y_hat_train_rp = lr_best_rand_precision.predict(X_train)
31 y_hat_test_rp = lr_best_rand_precision.predict(X_test)
```

Fitting 3 folds for each of 100 candidates, totalling 300 fits
231.75220894813538

```
-----
--
NameError                                Traceback (most recent call last)
<ipython-input-20-5d1ead3fffa2> in <module>
    32
    33
--> 34 print('Training Precision: ', precision_score(y_train_rp, y_hat_train_rp))
    35 print('Testing Precision: ', precision_score(y_test_rp, y_hat_test_rp))
    36 print('\n\n')

NameError: name 'y_train_rp' is not defined
```

```
In [21]: 1 print('Training Precision: ', precision_score(y_train, y_hat_train_rp))
2 print('Testing Precision: ', precision_score(y_test, y_hat_test_rp))
3 print('\n\n')
4
5
6 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_rp))
7 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_rp))
8 print('\n\n')
9
10 print('Training F1-Score: ', f1_score(y_train, y_hat_train_rp))
11 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_rp))
```

Training Precision: 0.9579011592434411

Testing Precision: 0.9559834938101788

Training Accuracy: 0.5780663780663781

Testing Accuracy: 0.5828843995510662

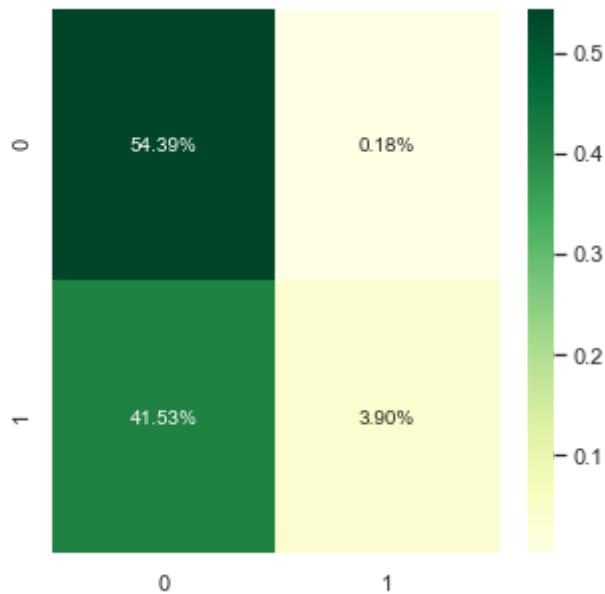
Training F1-Score: 0.15180816089731192

Testing F1-Score: 0.15754278590048737

```
In [22]: 1 #Get the confusion matrix
2 lr_rand_prec_matrix = confusion_matrix(y_test, y_hat_test_rp)
3 print(lr_rand_prec_matrix)
```

```
[[9692  32]
 [7401 695]]
```

```
In [55]: 1 # Visualize your confusion matrix
2 fig, ax = plt.subplots(figsize=(5,5))
3
4
5 sns.heatmap(lr_rand_prec_matrix/np.sum(lr_rand_prec_matrix), annot=True,
6             fmt='.2%', cmap='YlGn', ax=ax)
7
8 plt.show();
```



As the confusion matrix above makes it clear how poorly this model performs. Even though there is a very small amount of false positives, which is what we want, the model can hardly identify True Negatives at all! It misclassifies non-functional pumps at such a high rate that deployment of resources would go to very few water pumps that need them. This is very bad.

It is clear that focusing on just precision in order to lower false positives is a poor approach for building the model.

I will now focus on `f1_score` as my target for iterative modeling.

Below I decided to run a full grid search with the same parameters as my random grid, only this time with focus on improving `f1_score`.

```

In [25]: 1 start = time.time()
2
3 classifier_penalties = ['l1', 'l2']
4 classifier_Cs = np.logspace(-5, 5, 100)
5 classifier__solver = ['liblinear']
6
7 lr_grid = {'penalty' : classifier_penalties,
8           'C' : classifier_Cs,
9           'solver' : classifier__solver
10          }
11
12 lr = LogisticRegression()
13 # Random search of parameters, using 3 fold cross validation,
14 # search across 100 different combinations, and use all available cores
15 lr_grid_f1 = GridSearchCV(estimator=lr,
16                           param_grid=lr_grid,
17                           scoring='f1',
18                           cv=3,
19                           n_jobs = -1,
20                           verbose=2
21                          )
22
23 # RandomizedSearchCV(estimator = lr,
24 #                    param_distributions = random_grid,
25 #                    n_iter = 100, cv = 3, verbose=2,
26 #                    random_state=42,
27 #                    n_jobs = -1)
28 # Fit the random search model
29 lr_grid_f1.fit(X_train, y_train)
30
31
32 end = time.time()
33 print(end - start)
34
35
36 # lr_best_grid_f1 = lr_grid_f1.best_estimator_
37 # lr_best_grid_f1.fit(X_train,y_train)
38 # y_hat_train_gs = lr_best_grid_f1.predict(X_train)
39 # y_hat_test_gs = lr_best_grid_f1.predict(X_test)

```

Fitting 3 folds for each of 200 candidates, totalling 600 fits
469.99079990386963

```

-----
--
NameError                                Traceback (most recent call last)
<ipython-input-25-de0ab4e2fd38> in <module>
    35 lr_best_grid_f1 = lr_grid_f1.best_estimator_
    36 lr_best_grid_f1.fit(X_train,y_train)
--> 37 y_hat_train_gs = lr_best_rand_score.predict(X_train)
    38 y_hat_test_gs = lr_best_rand_score.predict(X_test)

NameError: name 'lr_best_rand_score' is not defined

```



```
In [29]: 1 lr_grid_f1.best_estimator_
```

```
Out[29]: LogisticRegression(C=5.72236765935022, solver='liblinear')
```

```
In [31]: 1 lr_best_grid_f1 = lr_grid_f1.best_estimator_  
2 lr_best_grid_f1.fit(X_train,y_train)  
3 y_hat_train_gs = lr_best_grid_f1.predict(X_train)  
4 y_hat_test_gs = lr_best_grid_f1.predict(X_test)  
5  
6  
7 print('Training Precision: ', precision_score(y_train, y_hat_train_gs))  
8 print('Testing Precision: ', precision_score(y_test, y_hat_test_gs))  
9 print('\n\n')  
10  
11  
12 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_gs))  
13 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_gs))  
14 print('\n\n')  
15  
16 print('Training F1-Score: ', f1_score(y_train, y_hat_train_gs))  
17 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_gs))
```

Training Precision: 0.7952357133925222

Testing Precision: 0.7879630985503002

Training Accuracy: 0.769095719095719

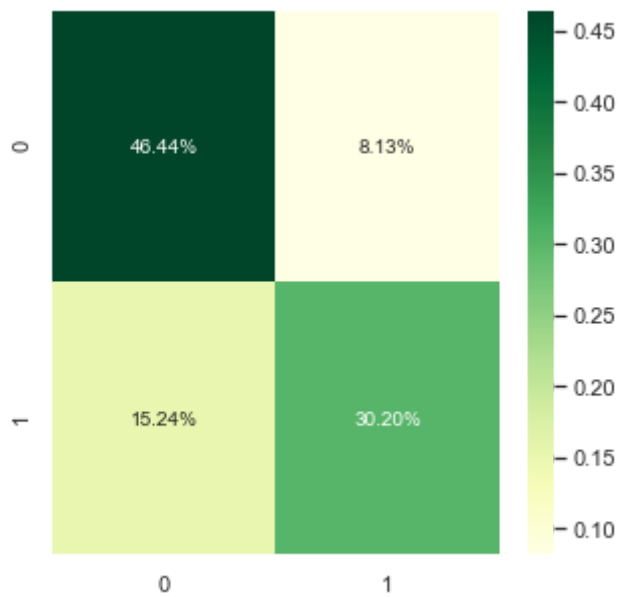
Testing Accuracy: 0.7663860830527497

Training F1-Score: 0.7259910385570364

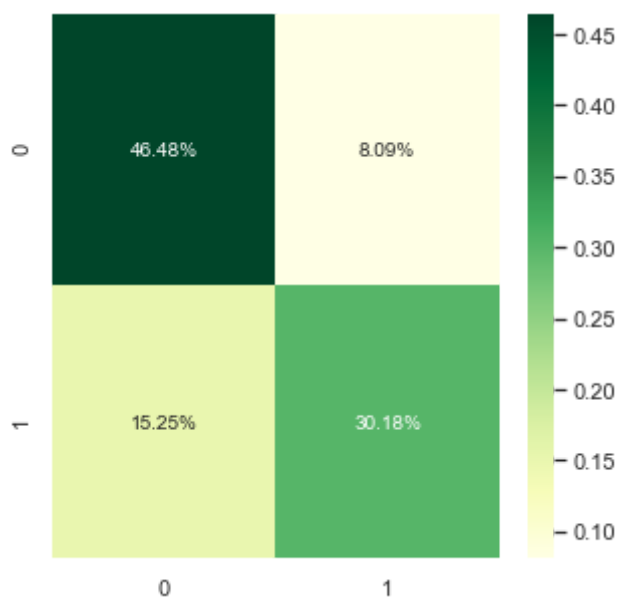
Testing F1-Score: 0.72107202680067

```
In [54]: 1 #best grid search confusion matrix
2 lr_grid_f1_matrix = confusion_matrix(y_test, y_hat_test_gs)
3 print(lr_grid_f1_matrix)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5,5))
7
8
9 sns.heatmap(lr_grid_f1_matrix/np.sum(lr_grid_f1_matrix), annot=True,
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();
```

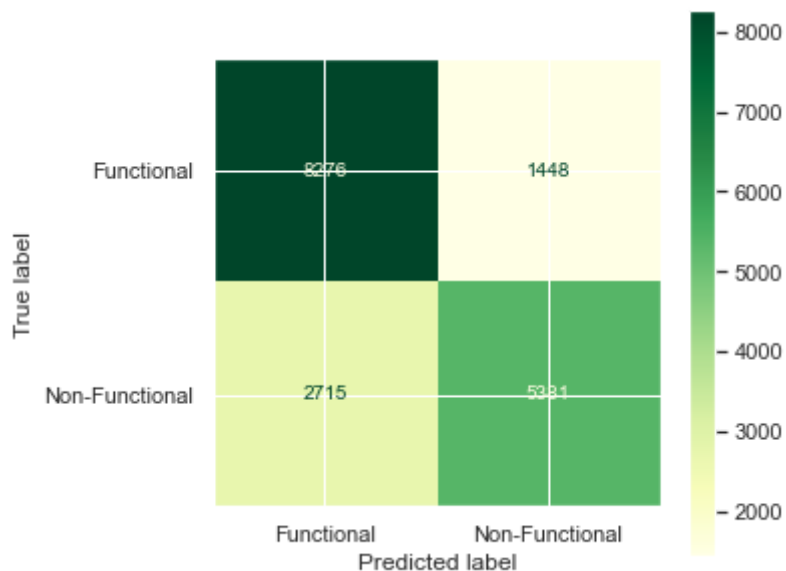
```
[[ 8276 1448]
 [ 2715 5381]]
```



```
In [52]: 1 #baseline confusion matrix
2 fig, ax = plt.subplots(figsize=(5,5))
3
4
5 sns.heatmap(lr_base_matrix/np.sum(lr_base_matrix), annot=True,
6             fmt='.2%', cmap='YlGn', ax=ax)
7
8 plt.show();
```



```
In [53]: 1 fig, ax = plt.subplots(figsize=(5,5))
2
3 plot_confusion_matrix(lr_best_grid_f1,
4                       X_test,
5                       y_test,
6                       ax=ax,
7                       cmap='YlGn',
8                       display_labels=fig_labels)
9
10 plt.show();
```



The gridsearch did not do much to improve the model.

Random Forest Classifier

Below I have created a grid of paramters in order to tune a random forest classifier. This grid will be used to perform a random grid search in order to narrow down the paramters to run a more robust grid search. The run time for a Random Forest has potential to be much higher than a logistic regression model. Therefore, a random search to narrow down the paramters for a grid search was carried out in order to reduce total run time while still keeping a fairly broad search.

```
In [40]: 1 # Number of trees in random forest
2 n_estimators = [int(x) for x in np.linspace(start = 200, stop = 2000, n
3 # Number of features to consider at every split
4 max_features = ['auto', 'sqrt']
5 # Maximum number of levels in tree
6 max_depth = [int(x) for x in np.linspace(10, 110, num = 11)]
7 max_depth.append(None)
8 # Minimum number of samples required to split a node
9 min_samples_split = [2, 5, 10]
10 # Minimum number of samples required at each leaf node
11 min_samples_leaf = [1, 2, 4]
12 # Create the random grid
13 random_grid = {'n_estimators': n_estimators,
14                 'max_features': max_features,
15                 'max_depth': max_depth,
16                 'min_samples_split': min_samples_split,
17                 'min_samples_leaf': min_samples_leaf}
18 random_grid
```

```
Out[40]: {'n_estimators': [200, 400, 600, 800, 1000, 1200, 1400, 1600, 1800, 2000],
'max_features': ['auto', 'sqrt'],
'max_depth': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, None],
'min_samples_split': [2, 5, 10],
'min_samples_leaf': [1, 2, 4]}
```

```

In [41]: 1 start = time.time()
2
3 # Use the random grid to search for best hyperparameters
4 # First create the base model to tune
5 rfc = RandomForestClassifier()
6 # Random search of parameters, using 3 fold cross validation,
7 # search across 100 different combinations, and use all available cores
8 rfc_random = RandomizedSearchCV(scoring='f1',
9                                 estimator = rfc,
10                                param_distributions = random_grid,
11                                n_iter = 33,
12                                cv = 3,
13                                verbose=2,
14                                random_state=42,
15                                n_jobs = -1)
16 # Fit the random search model
17
18 rfc_random.fit(X_train,y_train)
19 rfc_best_random = rfc_random.best_estimator_
20 rfc_best_random.fit(X_train,y_train)
21 y_hat_train_rrf = rfc_best_random.predict(X_train)
22 y_hat_test_rrf = rfc_best_random.predict(X_test)
23
24
25 end = time.time()
26 print(end - start)
27
28 print('Training Precision: ', precision_score(y_train, y_hat_train_rrf))
29 print('Testing Precision: ', precision_score(y_test, y_hat_test_rrf))
30 print('\n\n')
31
32
33 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_rrf))
34 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_rrf))
35 print('\n\n')
36
37 print('Training F1-Score: ', f1_score(y_train, y_hat_train_rrf))
38 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_rrf))

```

Fitting 3 folds for each of 33 candidates, totalling 99 fits

2072.605688095093

Training Precision: 0.9026457450463864

Testing Precision: 0.8049204439096136

Training Accuracy: 0.8801587301587301

Testing Accuracy: 0.8016273849607183

Training F1-Score: 0.863505629057441

Testing F1-Score: 0.7730337078651686

```
In [37]: 1 rfc_random.best_estimator_.get_params()
```

```
Out[37]: {'bootstrap': True,
'ccp_alpha': 0.0,
'class_weight': None,
'criterion': 'gini',
'max_depth': None,
'max_features': 'auto',
'max_leaf_nodes': None,
'max_samples': None,
'min_impurity_decrease': 0.0,
'min_impurity_split': None,
'min_samples_leaf': 1,
'min_samples_split': 10,
'min_weight_fraction_leaf': 0.0,
'n_estimators': 600,
'n_jobs': None,
'oob_score': False,
'random_state': None,
'verbose': 0,
'warm_start': False}
```

```
In [51]: 1 #best random search confusion matrix for Random Forest
2 lr_grid_f1_matrix_rrf = confusion_matrix(y_test, y_hat_test_rrf)
3 print(lr_grid_f1_matrix_rrf)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5,5))
7
8
9 sns.heatmap(lr_grid_f1_matrix_rrf/np.sum(lr_grid_f1_matrix_rrf), annot=
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();
```

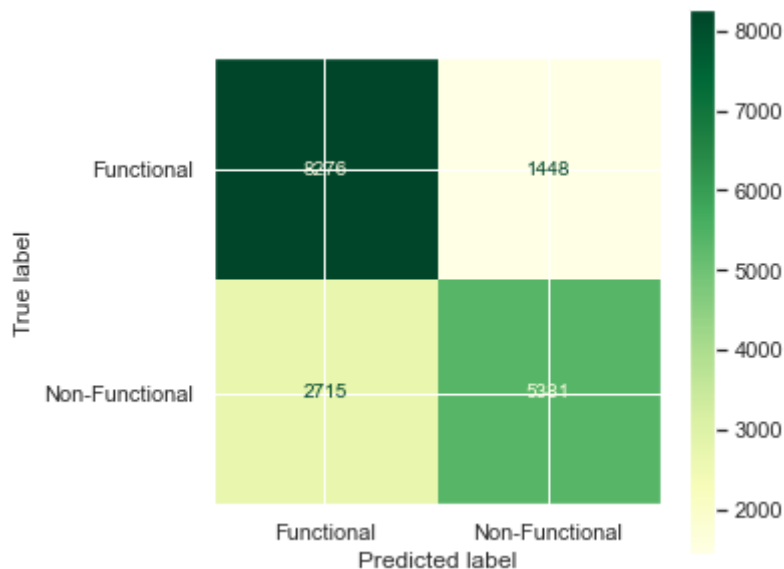
```
[[8265 1459]
 [2076 6020]]
```



```

In [50]: 1 fig, ax = plt.subplots(figsize=(5,5))
          2
          3 plot_confusion_matrix(lr_best_grid_f1,
          4                           X_test,
          5                           y_test,
          6                           ax=ax,
          7                           cmap='YlGn',
          8                           display_labels=fig_labels)
          9
          10 plt.show();

```



This random forest classifier outperforms the best logistic regression model in every relevant metric. It improves upon the model accuracy and precision while also lowering the occurrence false negatives.

A grid search will be performed below to see if any improvements can be made to model performance.

```

In [44]: 1 # Number of trees in random forest
          2 n_estimators = [500, 600, 700]
          3 # Number of features to consider at every split
          4 max_features = [14,15,16]
          5 # Maximum number of levels in tree
          6 # Minimum number of samples required to split a node
          7 min_samples_split = [8,10,12]
          8 # Minimum number of samples required at each leaf node
          9 # Create the random grid
          10 cv_grid = {'n_estimators': n_estimators,
          11                'max_features': max_features,
          12                'min_samples_split': min_samples_split}
          13
          14 cv_grid

```

```

Out[44]: {'n_estimators': [500, 600, 700],
          'max_features': [14, 15, 16],
          'min_samples_split': [8, 10, 12]}

```

```
In [45]: 1 start = time.time()
2
3 rfc = RandomForestClassifier()
4 # Instantiate the grid search model
5 rfc_gs = GridSearchCV(scoring='f1',
6                       estimator = rfc,
7                       param_grid = cv_grid,
8                       cv = 3,
9                       n_jobs = -1,
10                      verbose = 2)
11
12
13 rfc_gs.fit(X_train,y_train)
14 rfc_best_gs = rfc_gs.best_estimator_
15 rfc_best_gs.fit(X_train,y_train)
16 y_hat_train_grf = rfc_best_gs.predict(X_train)
17 y_hat_test_grf = rfc_best_gs.predict(X_test)
18
19
20 end = time.time()
21 print(end - start)
22
23 print('Training Precision: ', precision_score(y_train, y_hat_train_grf))
24 print('Testing Precision: ', precision_score(y_test, y_hat_test_grf))
25 print('\n\n')
26
27
28 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_grf))
29 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_grf))
30 print('\n\n')
31
32 print('Training F1-Score: ', f1_score(y_train, y_hat_train_grf))
33 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_grf))
```

Fitting 3 folds for each of 27 candidates, totalling 81 fits

868.361251115799

Training Precision: 0.9079278044322595

Testing Precision: 0.8040468583599574

Training Accuracy: 0.8854978354978355

Testing Accuracy: 0.802020202020202

Training F1-Score: 0.8697507728503817

Testing F1-Score: 0.7739620707329574


```

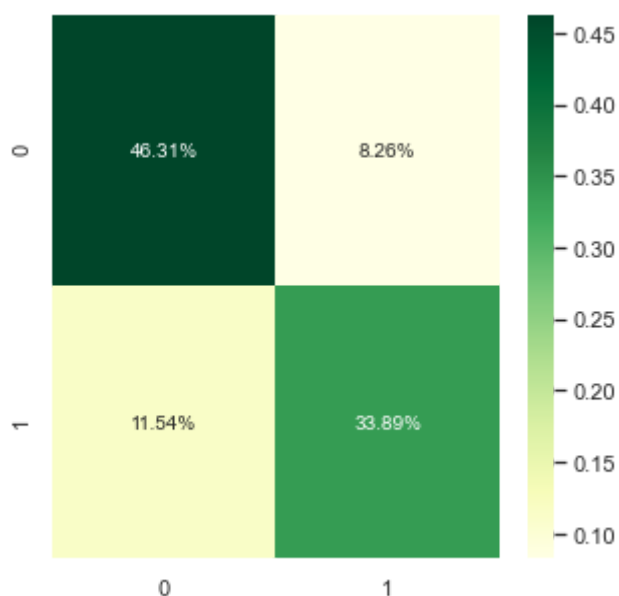
In [48]: 1 #best grid search confusion matrix for Random Forest
2 lr_grid_fl_matrix_grf = confusion_matrix(y_test, y_hat_test_grf)
3 print(lr_grid_fl_matrix_grf)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5, 5))
7
8
9 sns.heatmap(lr_grid_fl_matrix_grf/np.sum(lr_grid_fl_matrix_grf), annot=
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();
13
14 #best random search matrix
15 # [[8265 1459]
16 #  [2076 6020]]

```

```

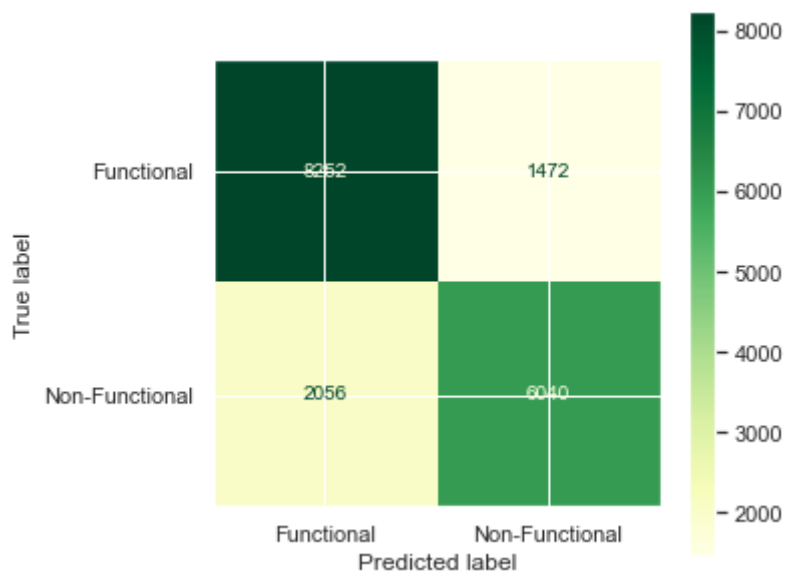
[[8252 1472]
 [2056 6040]]

```



The results of the grid search reduced the number of false positives, which is great, but it slightly increased the number of false negatives. It identified a similar number of functional pumps as the other random forest model, but it also identified slightly more non-functional pumps correctly. The modeling was slightly improved overall.

```
In [49]: 1 fig, ax = plt.subplots(figsize=(5, 5))
2
3 plot_confusion_matrix(rfc_best_gs,
4                       X_test,
5                       y_test,
6                       ax=ax,
7                       cmap='YlGn',
8                       display_labels=fig_labels)
9
10 plt.show();
```



XGBoost Classification Iterative Modeling

A similar method of iterative modeling was carried out using XGBoost as was used for tuning the random forest classifier: a random grid cv will be used to reduce the number of parameters to perform in a more robust grid search. A baseline XGBoost classifier was first created using default parameters before iterative modeling/hyper-parameter tuning was carried out.

```
In [58]: 1 # Instantiate the model
2 xgc_base = XGBClassifier()
3
4 # Fit the model
5 xgc_base.fit(X_train, y_train)
6 y_hat_train = xgc_base.predict(X_train)
7 y_hat_test = xgc_base.predict(X_test)
8
9 print('Training Precision: ', precision_score(y_train, y_hat_train))
10 print('Testing Precision: ', precision_score(y_test, y_hat_test))
11 print('\n\n')
12
13 print('Training Recall: ', recall_score(y_train, y_hat_train))
14 print('Testing Recall: ', recall_score(y_test, y_hat_test))
15 print('\n\n')
16
17 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train))
18 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test))
19 print('\n\n')
20
21 print('Training F1-Score: ', f1_score(y_train, y_hat_train))
22 print('Testing F1-Score: ', f1_score(y_test, y_hat_test))
```

Training Precision: 0.8397065362383281
Testing Precision: 0.8348687641230662

Training Recall: 0.5949593069046994
Testing Recall: 0.5932559288537549

Training Accuracy: 0.7624579124579125
Testing Accuracy: 0.7618967452300786

Training F1-Score: 0.696456559820523
Testing F1-Score: 0.6936240883818325

```

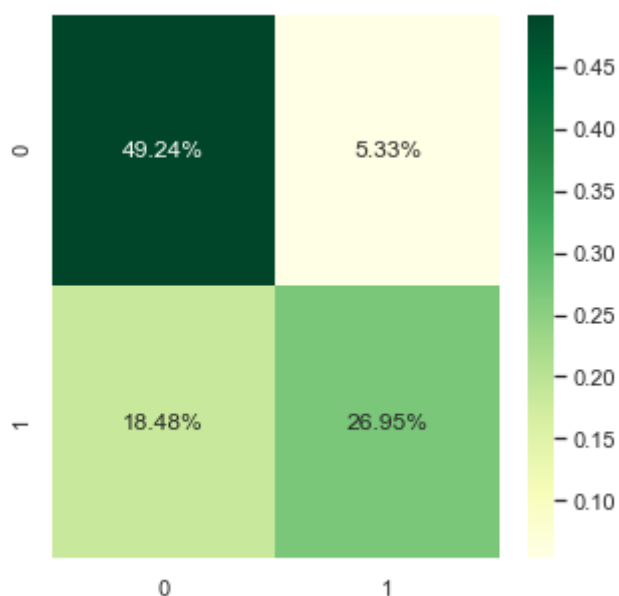
In [61]: 1 #baseline XGBoost model
2 xgs_base_matrix = confusion_matrix(y_test, y_hat_test)
3 print(xgs_base_matrix)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5, 5))
7
8
9 sns.heatmap(xgs_base_matrix/np.sum(xgs_base_matrix), annot=True,
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();

```

```

[[ 8774  950]
 [ 3293 4803]]

```



The baseline XGBoost model has a low number of false negatives, but a relatively high number of false positives. It also has a lower accuracy than either random forest model. Some hyper-parameter tuning will be carried out below in an attempt to improve the model.

```

In [62]: 1 params = {
2         'min_child_weight': [1, 5, 10],
3         'gamma': [0.5, 1, 1.5, 2, 5],
4         'subsample': [0.6, 0.8, 1.0],
5         'colsample_bytree': [0.6, 0.8, 1.0],
6         'max_depth': [3, 4, 5]
7     }
8
9     xgc = XGBClassifier()
10    # Random search of parameters, using 3 fold cross validation,
11    # search across 100 different combinations, and use all available cores
12    xgc_random = RandomizedSearchCV(scoring='f1',
13                                   estimator = xgc,
14                                   param_distributions = params,
15                                   n_iter = 33, cv = 3, verbose=2,
16                                   random_state=42,
17                                   n_jobs = -1)
18    # Fit the random search model
19    xgc_random.fit(X_train, y_train)
20
21
22    end = time.time()
23    print(end - start)
24
25
26    # xgc_best_random = xgc_random.best_estimator_
27    # xgc_best_random.fit(X_train,y_train)
28    # y_hat_train_brm = xgc_best_random.predict(X_train)
29    # y_hat_test_brm = xgc_best_random.predict(X_test)
30
31    # print('Training Precision: ', precision_score(y_train, y_hat_train_brm))
32    # print('Testing Precision: ', precision_score(y_test, y_hat_test_brm))
33    # print('\n\n')
34
35    # print('Training Recall: ', recall_score(y_train, y_hat_train_brm))
36    # print('Testing Recall: ', recall_score(y_test, y_hat_test_brm))
37    # print('\n\n')
38
39    # print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_brm))
40    # print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_brm))
41    # print('\n\n')
42
43    # print('Training F1-Score: ', f1_score(y_train, y_hat_train_brm))
44    # print('Testing F1-Score: ', f1_score(y_test, y_hat_test_brm))

```

Fitting 3 folds for each of 33 candidates, totalling 99 fits
 3434.718137025833

```

-----
--
NameError                                Traceback (most recent call last)
<ipython-input-62-89045d98e77d> in <module>
      26 xgc_best_random = xgc_random.best_estimator_
      27 xgc_best_random.fit(X_train,y_train)
--> 28 y_hat_train_brm = lr_best_rand_score.predict(X_train)

```

```

29 y_hat_test_brm = lr_best_rand_score.predict(x_test)
30

```

NameError: name 'lr_best_rand_score' is not defined

```

In [63]: 1 xgc_best_random = xgc_random.best_estimator_
          2 xgc_best_random.fit(X_train,y_train)
          3 y_hat_train_brm = xgc_best_random.predict(X_train)
          4 y_hat_test_brm = xgc_best_random.predict(X_test)
          5
          6 print('Training Precision: ', precision_score(y_train, y_hat_train_brm))
          7 print('Testing Precision: ', precision_score(y_test, y_hat_test_brm))
          8 print('\n\n')
          9
          10 print('Training Recall: ', recall_score(y_train, y_hat_train_brm))
          11 print('Testing Recall: ', recall_score(y_test, y_hat_test_brm))
          12 print('\n\n')
          13
          14 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_brm))
          15 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_brm))
          16 print('\n\n')
          17
          18 print('Training F1-Score: ', f1_score(y_train, y_hat_train_brm))
          19 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_brm))

```

Training Precision: 0.8566377816291161
 Testing Precision: 0.8387410772225827

Training Recall: 0.6488317143607246
 Testing Recall: 0.6385869565217391

Training Accuracy: 0.7894179894179895
 Testing Accuracy: 0.7800224466891134

Training F1-Score: 0.7383925903794443
 Testing F1-Score: 0.7251051893408134

```

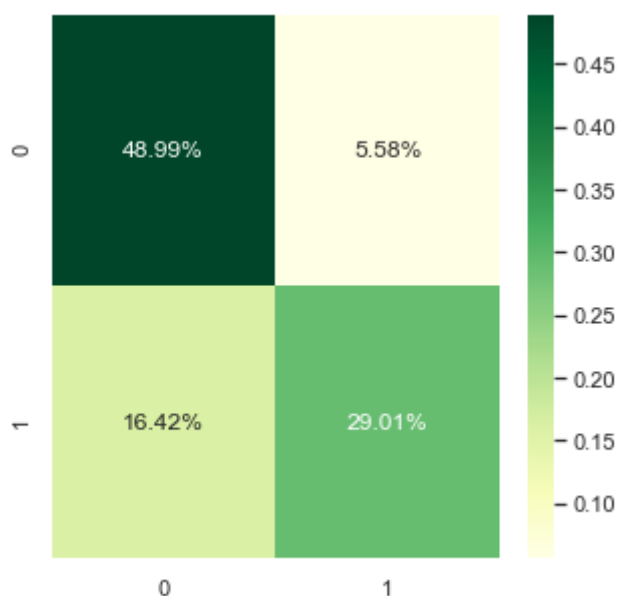
In [65]: 1 #baseline XGBoost model
2 xgs_rand_matrix = confusion_matrix(y_test, y_hat_test_brm)
3 print(xgs_rand_matrix)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5, 5))
7
8
9 sns.heatmap(xgs_rand_matrix/np.sum(xgs_rand_matrix), annot=True,
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();

```

```

[[8730  994]
 [2926 5170]]

```



Improvements to accuracy and f1-score. Still needs improvement to both to compete with the random forest models.

```

In [64]: 1 xgc_best_random

```

```

Out[64]: XGBClassifier(colsample_bytree=0.8, gamma=0.5, max_depth=5, subsample=1.0)

```

```

In [66]: 1 params2 = {
2         'gamma': [0.1, 0.3, 0.5],
3         'colsample_bytree': [0.8, 0.9],
4         'max_depth': [4, 5, 6]
5         }
6
7 start = time.time()
8
9 xgc2 = XGBClassifier()
10 # Random search of parameters, using 3 fold cross validation,
11 # search across 100 different combinations, and use all available cores
12 xgc_gs = GridSearchCV(scoring='f1',
13                       estimator = xgc2,
14                       param_grid = params2,
15                       cv = 3,
16                       verbose=2,
17                       n_jobs = -1)
18
19
20 xgc_gs.fit(X_train, y_train)
21 xgc_best_gs = xgc_gs.best_estimator_
22 xgc_best_gs.fit(X_train, y_train)
23 y_hat_train_bgs = xgc_best_gs.predict(X_train)
24 y_hat_test_bgs = xgc_best_gs.predict(X_test)
25
26
27 end = time.time()
28 print(end - start)
29
30 print('Training Precision: ', precision_score(y_train, y_hat_train_bgs))
31 print('Testing Precision: ', precision_score(y_test, y_hat_test_bgs))
32 print('\n\n')
33
34
35 print('Training Accuracy: ', accuracy_score(y_train, y_hat_train_bgs))
36 print('Testing Accuracy: ', accuracy_score(y_test, y_hat_test_bgs))
37 print('\n\n')
38
39 print('Training F1-Score: ', f1_score(y_train, y_hat_train_bgs))
40 print('Testing F1-Score: ', f1_score(y_test, y_hat_test_bgs))

```

Fitting 3 folds for each of 18 candidates, totalling 54 fits

553.8524680137634

Training Precision: 0.864413680781759

Testing Precision: 0.8411867364746946

Training Accuracy: 0.8002645502645502

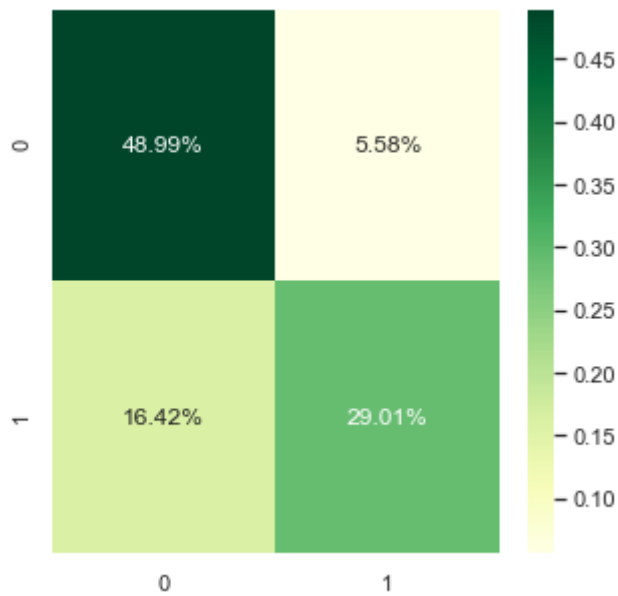
Testing Accuracy: 0.7870370370370371

Training F1-Score: 0.7541517421035494

Testing F1-Score: 0.7364400305576775


```
In [67]: 1 #baseline XGBoost model
2 xgs_rand_matrix = confusion_matrix(y_test, y_hat_test_brm)
3 print(xgs_rand_matrix)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5, 5))
7
8
9 sns.heatmap(xgs_rand_matrix/np.sum(xgs_rand_matrix), annot=True,
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();
```

```
[[ 8730  994]
 [ 2926 5170]]
```



```

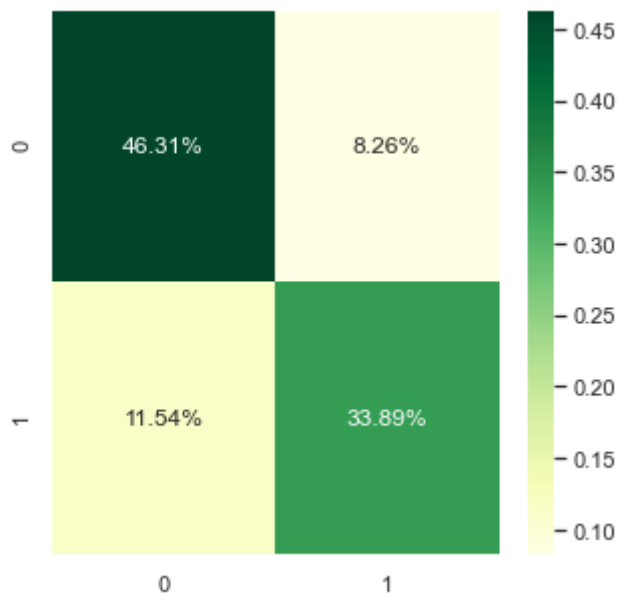
In [68]: 1 #best grid search confusion matrix for Random Forest
2 lr_grid_fl_matrix_grf = confusion_matrix(y_test, y_hat_test_grf)
3 print(lr_grid_fl_matrix_grf)
4
5 # Visualize your confusion matrix
6 fig, ax = plt.subplots(figsize=(5, 5))
7
8
9 sns.heatmap(lr_grid_fl_matrix_grf/np.sum(lr_grid_fl_matrix_grf), annot=
10             fmt='.2%', cmap='YlGn', ax=ax)
11
12 plt.show();
13
14 #best random search matrix
15 # [[8265 1459]
16 #  [2076 6020]]

```

```

[[8252 1472]
 [2056 6040]]

```



The grid search produced a pretty good model, but my random forest classifier still performed better.

Evaluation

My random forest models outperformed my best logistic regression and XGBoost models in regards to the metrics that are most important given the business problem at hand.

The best random forest model has the lowest false positive rate, a low false negative rate, and the highest accuracy.

About 11.5% of pumps would be misclassified as functional using my best model. This means that 11.5% of the pumps would go untreated if it was deployed to conduct predictive maintenance. However, it correctly identifies a high number of functional pumps correctly, which would save a lot

of valuable resources, time and money, and it also identifies a large number of non-functional correctly. Only about 8.25% of functional pumps would be incorrectly identified as non-functional. This is the resource/time/money sink of my model.

Conclusions

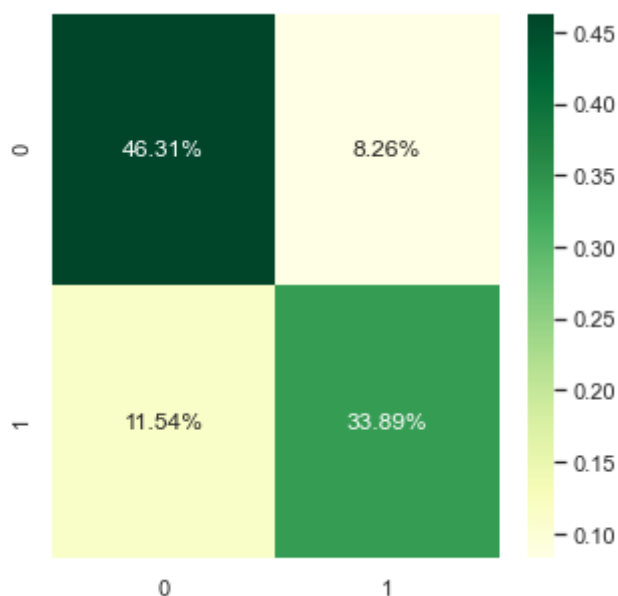
I believe that my best classification model provides a powerful enough predictive ability to prove very valuable to the Ministry of Water. The amount of resources saved, the relatively low number of missclassified functional pumps, and the elimination of the need to physically sweep the functionality of all pumps can bring access to potable drinking water to a larger number of communities than before without predictive maintenance.

```
In [76]: 1 df_binary['status_group'].value_counts(normalize=True)
```

```
Out[76]: 0    0.543081
         1    0.456919
         Name: status_group, dtype: float64
```

```
In [72]: 1 #best grid search confusion matrix for Random Forest
         2 print(lr_grid_fl_matrix_grf)
         3
         4 # Visualize your confusion matrix
         5 fig, ax = plt.subplots(figsize=(5, 5))
         6
         7
         8 sns.heatmap(lr_grid_fl_matrix_grf/np.sum(lr_grid_fl_matrix_grf), annot=
         9             fmt='.2%', cmap='YlGn', ax=ax)
        10
        11 plt.show();
```

```
[[ 8252 1472]
 [ 2056 6040]]
```



Thank you! For questions or comments please feel free to reach me by email.

Author: Dylan Dey

Email: ddey2985@gmail.com (<mailto:ddey2985@gmail.com>).