

PREDICTING THE CONDITION OF WATER PUMPS IN TANZANIA



Business Understanding

The Problem

Tanzania is a developing country in eastern Africa. It is a geographically diverse country with mountainous terrain and flat plains. The country borders the Indian Ocean in the east, and the Great Rift Valley on its western border. A good portion of the country is below sea level, and much of it is 900 ft above sea level. Access to natural resources is spotty, and water is one of the most, if not the most, scarce resource in the country. Many organizations have installed water pumps in villages around the country in an effort to provide clean drinking water to the people. These pumps vary in how they extract water, the water quality, and what basin the water comes from, how they are managed, the population around the well, and whether the users have to pay for water or not. Pumps break down regularly and require maintenance which can be quite a task considering there are over 59,000 pumps in the country.

The Task

We have been asked by the Tanzanian Government to create a prediction model that will predict the condition of a water pump; functional, or non-functional. This will allow them to strategically mobilize repair teams in areas with a high concentration of non-functional pumps, and to efficiently

react when resources are needed in a particular area, such as supplying potable water to a village whose pump is non-functional and awaiting repair.

Hypothesis

Our null hypothesis is that we CANNOT predict whether a well is functional or not.

Our alternative hypothesis is that we CAN predict the condition of a well.

- A false positive would be to predict a well is non-functional when it is functional.
- A false negative would be to predict a well is functional when it is not functional.

Data Understanding

The dataset is provided by Taarifa which aggregates data from the Tanzania Ministry of Water on the over 59,000 water pumps in the country. The data contains information about the location, operation, management, installation, water quality, and population of users of a particular well.

Method:

- Determine relevant features for prediction
- Create several models and look for best predictor metrics
- Use validation methods to ensure performance
- Tune final model for optimal predictions

Model Selection

NOTE: a positive prediction is a well that is non-functional.

Maximize Recall

We are seeking a model that will primarily maximize Recall. Recall will measure how well we are predicting wells that are actually non-functional. It is calculated by dividing the number of non-functioning wells that were correctly predicted by the total number of non-functional wells in our test set. Recall does not account for false positives.

Recall:

- Minimize false negatives
- Resources are directed to the people who need them
- More human lives are saved

With human lives at stake and counting on our model to get it right we will place primary importance on the recall score.

Maximize Precision

Secondarily we will be seeking a model that will maximize precision, but not at the expense of recall. Precision is a measure of how accurate the positive predictions are. It is calculated by dividing the number of correctly predicted positives by the total number of positive predictions. Precision does not account for false negatives.

Precision:

- Minimize false positives
- Resources are not directed where they are not needed
- Less logistic strain on the system, resources and man power

Repair materials and man power are also a scarce resource in this problem and we cannot afford to be sending repair crews out to wells that have been predicted as non-functional but are functional. However, human life is tantamount and we will not sacrifice recall for more precision.

The two metrics have an inverse relationship and it will be tricky to design a model that will maximize both.

Explore the data and process it for modeling

```
In [1]: 1 # IMPORT DEPENDENCIES
2
3 import numpy as np
4 import pandas as pd
5 import missingno as msno
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import yellowbrick as yb
9 import folium
10
11 from scipy import stats as stats
12 from sklearn.linear_model import LogisticRegression
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.tree import DecisionTreeClassifier
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.metrics import confusion_matrix, plot_confusion_matrix,\
17     precision_score, recall_score, accuracy_score, f1_score,\
18     classification_report
19 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
20 from sklearn.preprocessing import OneHotEncoder, StandardScaler
21 from sklearn.naive_bayes import GaussianNB
22
23 sns.set(style="whitegrid")
24 pd.set_option('display.max_columns', None)
25 pd.set_option("mode.chained_assignment", None)
```

```
In [2]: 1 # IMPORT THE DATA
        2
        3 V = pd.read_csv('Data/values.csv')
        4 y = pd.read_csv('Data/labels.csv')
```

```
In [3]: 1 # TRANSFORM OUR LABELS INTO A BINARY CLASS
        2
        3 y = y.replace({'status_group': {'functional' : 0, 'non functional' : 1,
```

```
In [4]: 1 # LOOK FOR IMBALANCE IN THE DATA
        2 # This looks good and should not cause an issue, hwever we will use bal
        3
        4 y.status_group.value_counts()
```

```
Out[4]: 0    32259
        1    27141
        Name: status_group, dtype: int64
```

```
In [5]: 1 # USE THE LATITUDE AND LONGITUDE AND COMBINE IT WITH OUR LABELS TO PLOT
        2
        3 df_lat_long = V['latitude'].to_frame().join(V['longitude']).join(y)
        4 map_center = [df_lat_long['latitude'].mean(), df_lat_long['longitude'].
```

```
In [45]: 1 # PLOT THE MAP
        2 # Notice the opacity denotes concentration of wells
        3
        4 map1 = folium.Map(location = map_center, tiles='Openstreetmap', zoom_st
        5
        6 for index, loc in df_lat_long.iterrows():
        7     if loc['status_group']==0:
        8         color = 'green'
        9     elif loc['status_group']==1:
        10        color = "red"
        11    else:
        12        color = 'white'
        13
        14        folium.CircleMarker([loc['latitude'], loc['longitude']], radius=2,
        15
        16 folium.LayerControl().add_to(map1)
        17
        18 map1
```

Out[45]: Make this Notebook Trusted to load map: File -> Trust Notebook

```
In [7]: 1 # Create a dataframe with our of our data
        2
        3 df_all = V.join(y.set_index('id'), on='id')
```

```
In [8]: 1 # Explore all of our data in one place
        2
        3 df_all.head()
```

```
Out[8]:
```

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

```
In [9]: 1 # Drop an outlier we found on the map (zoom out)
        2
        3 df_all.drop(labels=df_all.id[49651], inplace=True)
```

```
In [10]: 1 # No duplicates
         2
         3 df_all.duplicated().value_counts()
```

```
Out[10]: False      59399
dtype: int64
```

```
In [11]: 1 # Drop features that have no influence, and are repetitive
         2
         3 df = df_all.drop(columns=['id', 'amount_tsh', 'date_recorded', 'funder',
         4                          'longitude', 'latitude', 'wpt_name', 'num_private',
         5                          'subvillage', 'region', 'region_code', 'district_code',
         6                          'ward', 'public_meeting', 'recorded_by',
         7                          'scheme_management', 'scheme_name', 'permit',
         8                          'extraction_type_group', 'extraction_type_class',
         9                          'management_group', 'payment_type',
        10                          'quality_group', 'quantity', 'quantity_group',
        11                          'source_type', 'source_class',
        12                          'waterpoint_type_group'])
```

```
In [12]: 1 df.columns
```

```
Out[12]: Index(['gps_height', 'installer', 'basin', 'lga', 'population',  
              'construction_year', 'extraction_type', 'management', 'payment',  
              'water_quality', 'source', 'waterpoint_type', 'status_group'],  
              dtype='object')
```

```
In [13]: 1 # Check again for duplicates  
        2  
        3 df.duplicated().value_counts()
```

```
Out[13]: False      44034  
        True       15365  
        dtype: int64
```

We have duplicates now that we've dropped the 'id' column, and need to make a decision whether or not to drop them. We will not because we know they are independent wells with different 'id' numbers and they need to be represented when we train our models. It seems like we are overfitting models with repetitive information, to drop them would effectively be under-representing certain features and our models will not be accurate.

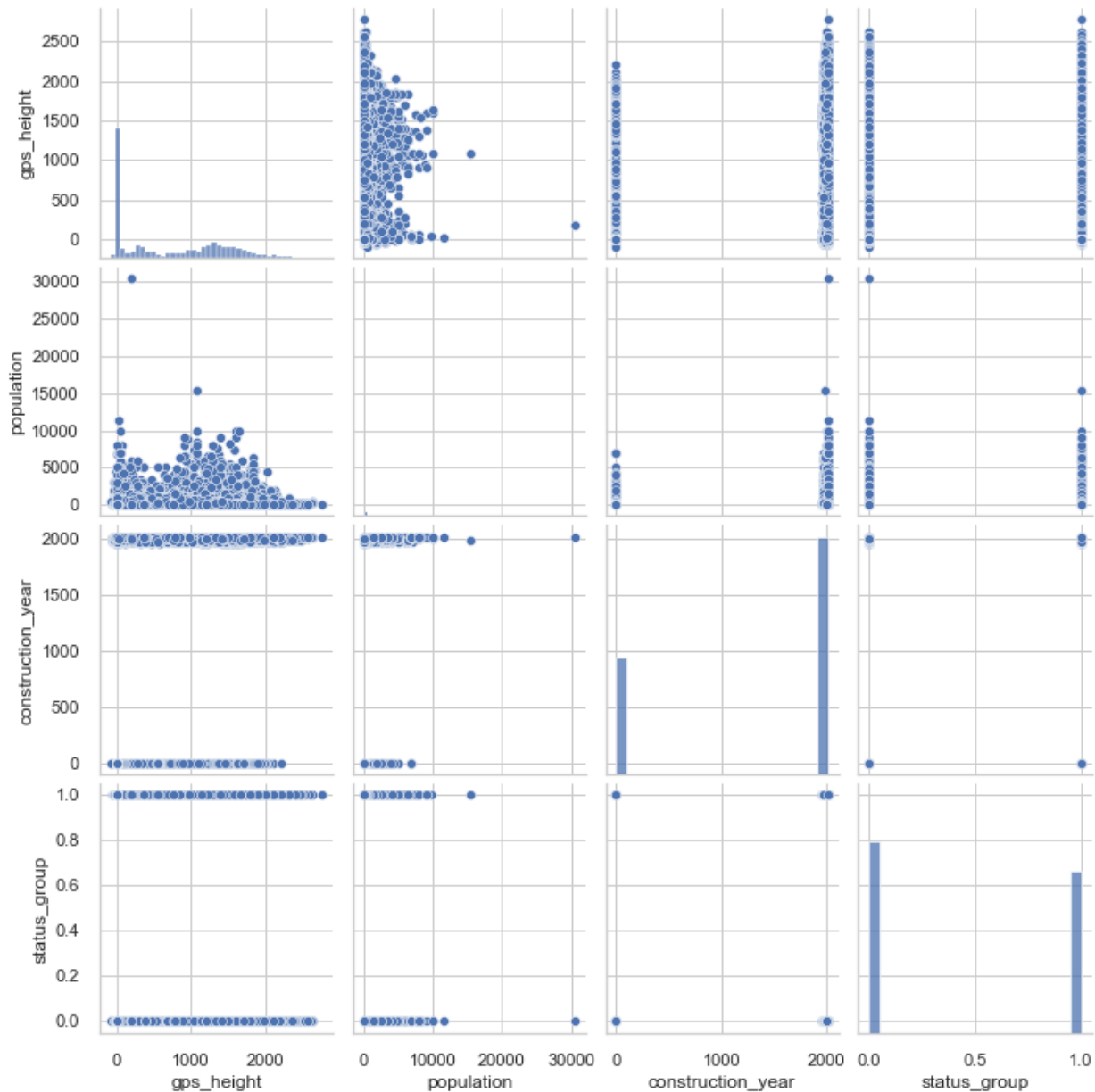
```
In [14]: 1 # df.drop_duplicates(inplace=True)
```

```
In [15]: 1 # Look for Null values  
        2  
        3 df.isna().sum().sum()
```

```
Out[15]: 3655
```

```
In [16]: 1 # Drop the Null values  
        2  
        3 df.dropna(inplace=True)
```

```
In [17]: 1 # Plot our numerical values against each other to look for any patterns
2
3 sns.pairplot(df)
4 plt.show()
```



Not much can be drawn about patterns, but we find some interesting information, however there are many 'construction_year' and 'population' entries with 0 values. 'gps_height' also shows many 0 values and even negative values, however this is in line with what we know about Tanzania. Some areas of the country are in valleys below sea level, and much of the country is above 900 ft. the median of our data set is approx. 990.

We will fix the construction year entries by imputing the median year, but the population 0 values may be accurate and reflect the population around a well and tell a story about the data, like perhaps an area has been abandoned due to lack of water, etc...

```
In [18]: 1 # Create our feature and label datasets for modeling
          2
          3 x = df.drop('status_group', axis=1)
          4 y = df.status_group
```



```

In [19]: 1  # Create a function that will process our data and split it into train
2
3  # A function that will process and scale numerical data
4  def process_scale (X, y):
5      # Encode categorical data
6      X_cat = X.select_dtypes('object')
7      ohe = OneHotEncoder(sparse=False, handle_unknown='ignore')
8      dums = ohe.fit_transform(X_cat)
9      dums_df = pd.DataFrame(dums, columns=ohe.get_feature_names(), index
10
11      # combine encoded columns back with numerical columns
12      X_nums = X.select_dtypes('int64')
13      X = pd.concat([X_nums, dums_df], axis=1)
14
15      # Split the data into training and test sets, we will test on only
16      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
17
18      # We do the following processing on train and test sets separately
19
20      # Replace 0's in 'construction_year' with the column median
21      X_train.construction_year.where(X_train.construction_year != 0, X_t
22
23      # Scale numerical data and combine back with encoded categorical co
24      X_train_nums = X_train.select_dtypes('int64')
25      ss = StandardScaler()
26      ss.fit(X_train_nums)
27      nums_df = pd.DataFrame(ss.transform(X_train_nums),
28                             index=X_train_nums.index)
29      X_cats = X_train.select_dtypes('float64')
30      X_train_clean = pd.concat([nums_df, X_cats], axis=1)
31
32      # Replace 0's in 'construction_year' with the column median
33      X_test.construction_year.where(X_test.construction_year != 0, X_tes
34
35      # Scale numerical data and combine back with encoded categorical co
36      X_test_nums = X_test.select_dtypes('int64')
37      ss = StandardScaler()
38      ss.fit(X_test_nums)
39      test_nums_df = pd.DataFrame(ss.transform(X_test_nums),
40                                  index=X_test_nums.index)
41      X_test_cats = X_test.select_dtypes('float64')
42      X_test_clean = pd.concat([test_nums_df, X_test_cats], axis=1)
43
44      return X_train_clean, X_test_clean, y_train, y_test
45
46
47  # The same process as above but with no scaling for tree-based models t
48  def process_no_scale (X, y):
49      X_cat = X.select_dtypes('object')
50      ohe = OneHotEncoder(sparse=False, handle_unknown='ignore')
51      dums = ohe.fit_transform(X_cat)
52      dums_df = pd.DataFrame(dums, columns=ohe.get_feature_names(), index
53
54      X_nums = X.select_dtypes('int64')
55      X = pd.concat([X_nums, dums_df], axis=1)
56

```

```

57     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
58
59     X_train.construction_year.where(X_train.construction_year != 0, X_t
60
61     X_test.construction_year.where(X_test.construction_year != 0, X_tes
62
63     return X_train, X_test, y_train, y_test

```

```

In [20]: 1 # Load our training and testing sets into memory
          2
          3 X_train, X_test, y_train, y_test = process_scale(X, y)
          4
          5 X_train_ns, X_test_ns, y_train_ns, y_test_ns = process_no_scale(X, y)

```

CREATE SOME MODELS AND SCORE THEM

We are looking to maximize Recall primarily.

We are also looking to maximize Precision, but not at the expense of Recall

Logistic Regression

```

In [21]: 1 # Instantiate and fit the model
          2
          3 logreg = LogisticRegression(fit_intercept=False, class_weight='balanced
          4 model_log = logreg.fit(X_train, y_train)

```

```

In [22]: 1 print(f"training accuracy: {model_log.score(X_train, y_train)}")

training accuracy: 0.7462177838904502

```

```

In [23]: 1 # Create predictions and score them
          2
          3 y_hat_test = model_log.predict(X_test)
          4
          5 print(confusion_matrix(y_test, y_hat_test))

```

```

[[2310  675]
 [ 749 1841]]

```

```
In [24]: 1 precision = precision_score(y_test, y_hat_test)
2 recall = recall_score(y_test, y_hat_test)
3 accuracy = accuracy_score(y_test, y_hat_test)
4 F1 = f1_score(y_test, y_hat_test)
5
6 print(precision)
7 print(accuracy)
8 print(recall)
9 print(F1)
```

```
0.7317170111287759
0.7445739910313901
0.7108108108108108
0.7211124167645906
```

K Nearest Neighbors

```
In [25]: 1 # Instantiate and fit the model
2
3 knn = KNeighborsClassifier()
4 knn.fit(X_train, y_train)
```

```
Out[25]: KNeighborsClassifier()
```

```
In [26]: 1 print(f"training accuracy: {knn.score(X_train, y_train)}")
```

```
training accuracy: 0.8259283621359804
```

```
In [27]: 1 # Create predictions and score them
2
3 knn_preds = knn.predict(X_test)
4
5 print(confusion_matrix(y_test, knn_preds))
```

```
[[2449  536]
 [ 675 1915]]
```

```
In [28]: 1 precision = precision_score(y_test, knn_preds)
2 recall = recall_score(y_test, knn_preds)
3 accuracy = accuracy_score(y_test, knn_preds)
4 F1 = f1_score(y_test, knn_preds)
5
6 print(precision)
7 print(accuracy)
8 print(recall)
9 print(F1)
```

```
0.7813137494900041
0.782780269058296
0.7393822393822393
0.7597698869271969
```

Those are some good scores

Recall: 74%
Precision: 78%

Naive Bayes

```
In [29]: 1 # Instantiate and fit the model
          2
          3 gnb = GaussianNB()
          4
          5 gnb.fit(X_train_ns, y_train)
          6
          7 print(f"training accuracy: {gnb.score(X_train, y_train)}")
          8 # print(f"testing accuracy: {gnb.score(X_test, y_test)}")
```

training accuracy: 0.5396360302178637

```
In [30]: 1 # Create predictions and score them
          2
          3 gnb_preds = gnb.predict(X_test)
          4 print(confusion_matrix(y_test, gnb_preds))
```

```
[[ 476 2509]
 [   71 2519]]
```

```
In [31]: 1 precision = precision_score(y_test, gnb_preds)
          2 recall = recall_score(y_test, gnb_preds)
          3 accuracy = accuracy_score(y_test, gnb_preds)
          4 F1 = f1_score(y_test, gnb_preds)
          5
          6 print(precision)
          7 print(accuracy)
          8 print(recall)
          9 print(F1)
```

0.500994431185362
0.537219730941704
0.9725868725868726
0.6613284326594907

Decision Tree

```
In [32]: 1 # Instantiate and fit the model
2
3 tree = DecisionTreeClassifier(max_depth=10, criterion='gini')
4
5 tree.fit(X_train_ns, y_train_ns)
6
7 print(f"training accuracy: {tree.score(X_train_ns, y_train_ns)}")
8 # print(f"testing accuracy: {tree.score(X_test_ns, y_test_ns)}")
```

training accuracy: 0.7315473698897725

```
In [33]: 1 # Create predictions and score them
2
3 tree_preds = tree.predict(X_test_ns)
4 print(confusion_matrix(y_test_ns, tree_preds))
```

```
[[2346  639]
 [ 938 1652]]
```

```
In [34]: 1 precision = precision_score(y_test_ns, tree_preds)
2 recall = recall_score(y_test_ns, tree_preds)
3 accuracy = accuracy_score(y_test_ns, tree_preds)
4 F1 = f1_score(y_test_ns, tree_preds)
5
6 print(precision)
7 print(accuracy)
8 print(recall)
9 print(F1)
```

```
0.7210824967263204
0.7171300448430493
0.6378378378378379
0.6769104691661545
```

Random Forest

```
In [35]: 1 # Instantiate and fit the model
2
3 forest = RandomForestClassifier(max_features=None, class_weight="balanced",
4                                max_depth = 10)
5
6 forest.fit(X_train_ns, y_train_ns)
7
8 print(f"training accuracy: {forest.score(X_train_ns, y_train_ns)}")
9 # print(f"testing accuracy: {forest.score(X_test_ns, y_test_ns)}")
```

training accuracy: 0.7529151468038031

```
In [36]: 1 # Create predictions and score them
          2
          3 forest_preds = forest.predict(X_test_ns)
          4 print(confusion_matrix(y_test_ns, forest_preds))
```

```
[[2429  556]
 [ 888 1702]]
```

```
In [37]: 1 precision = precision_score(y_test_ns, forest_preds)
          2 recall = recall_score(y_test_ns, forest_preds)
          3 accuracy = accuracy_score(y_test_ns, forest_preds)
          4 F1 = f1_score(y_test_ns, forest_preds)
          5
          6 print(precision)
          7 print(accuracy)
          8 print(recall)
          9 print(F1)
```

```
0.7537643932683791
0.7409865470852018
0.6571428571428571
0.7021452145214522
```

Based on Recall and Precision scores, our best performing model was the *K Nearest Neighbor model*

Tuning the KNN Model

We will iterate through a couple grid searches to be sure we are finding the best parameters

```
In [38]: 1 knn_tuned = KNeighborsClassifier(algorithm='auto')
2
3 parameters_KNN = {
4     'n_neighbors': (5, 23, 33, 43, 53),
5     'weights': ('uniform', 'distance')
6 }
7
8 grid_search_KNN = GridSearchCV(
9     estimator=knn_tuned,
10    param_grid=parameters_KNN,
11    scoring = 'recall',
12    n_jobs = -1,
13    cv = 5
14 )
15
16 KNN = grid_search_KNN.fit(X_train, y_train)
17
18 print('Best parameters:', grid_search_KNN.best_params_)
19 print('-----')
20 print('Best Recall Score:', grid_search_KNN.best_score_)
```

Best parameters: {'n_neighbors': 5, 'weights': 'distance'}

Best Recall Score - KNN: 0.7173736947961599

Lets fine tune this with n-neighbors closer to the above result

```
In [39]: 1 knn_tuned = KNeighborsClassifier(algorithm='auto')
2
3 parameters_KNN = {
4     'n_neighbors': (1, 3, 5, 7, 9, 11, 13),
5     'weights': ('uniform', 'distance')
6 }
7
8 grid_search_KNN = GridSearchCV(
9     estimator=knn_tuned,
10    param_grid=parameters_KNN,
11    scoring = 'recall',
12    n_jobs = -1,
13    cv = 5
14 )
15
16 KNN = grid_search_KNN.fit(X_train, y_train)
17
18 print('Best parameters:', grid_search_KNN.best_params_)
19 print('-----')
20 print('Best Recall Score:', grid_search_KNN.best_score_)
```

Best parameters: {'n_neighbors': 5, 'weights': 'distance'}

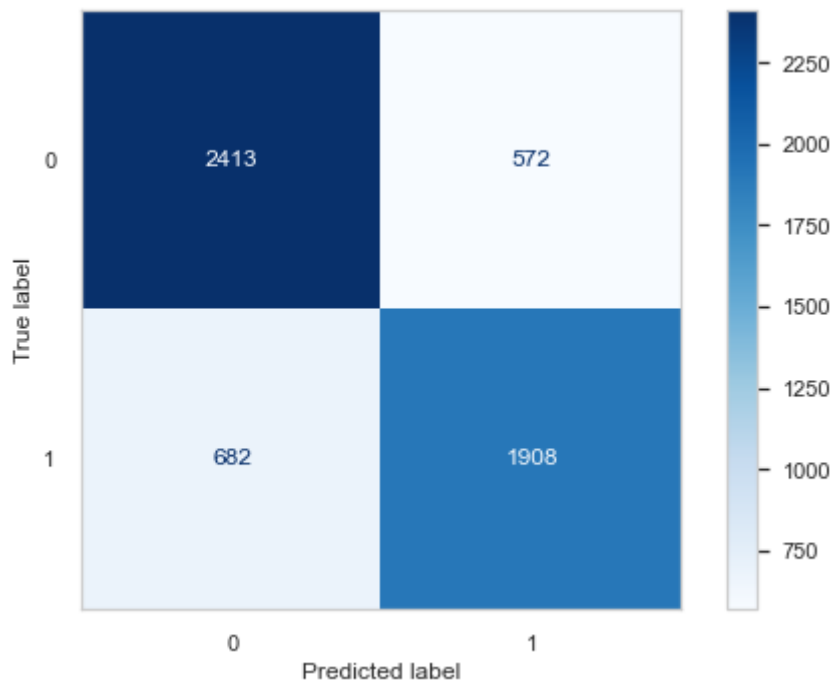
Best Recall Score - KNN: 0.7173736947961599

It looks like our default number of neighbors (5) was the top scoring number parameter, along with being weighted by distance. Let's see

how these parameters score on the test set.

distance weighting means that within "n" number of neighbors, closer neighbors are weighted heavier than neighbors that are further away

```
In [40]: 1 knn_tuned = KNeighborsClassifier(n_neighbors=5, weights='distance')
          2
          3 knn_tuned.fit(X_train, y_train)
          4
          5 knn_tuned_preds = knn_tuned.predict(X_test)
          6
          7 plot_confusion_matrix(knn_tuned, X_test, y_test, cmap='Blues')
          8 plt.grid(False)
```



```
In [41]: 1 precision = precision_score(y_test, knn_tuned_preds)
          2 recall = recall_score(y_test, knn_tuned_preds)
          3
          4 print('Precision:', precision)
          5 print('Recall:', recall)
```

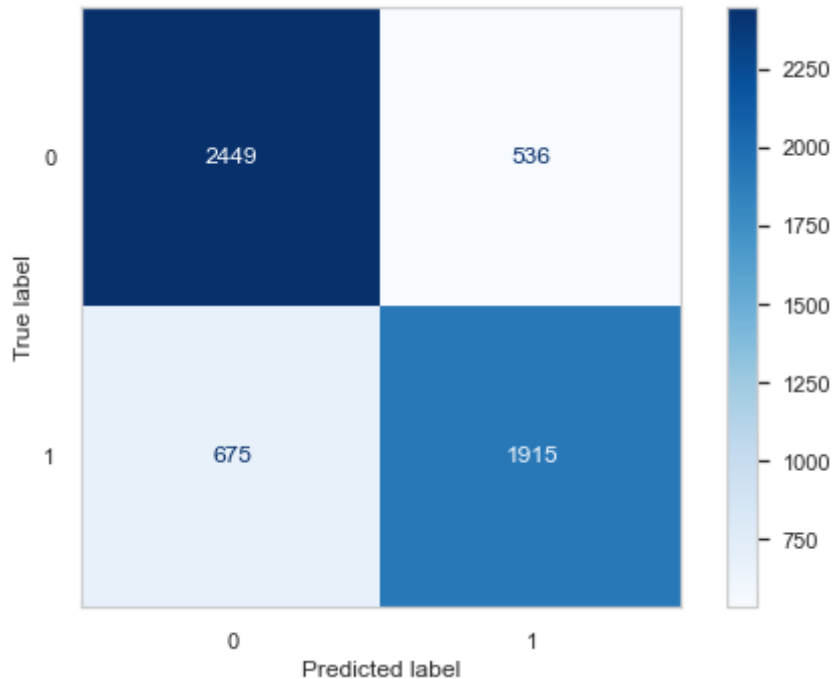
Precision: 0.7693548387096775
Recall: 0.7366795366795367

Let's compare these scores with our original default test scores.

Only difference is that uniform weighting is the default

Uniform means all neighbors within "n" are weighted equally


```
In [42]: 1 knn_tuned = KNeighborsClassifier(n_neighbors=5, weights='uniform')
2 knn_tuned.fit(X_train, y_train)
3 knn_tuned_preds = knn_tuned.predict(X_test)
4
5 plot_confusion_matrix(knn_tuned, X_test, y_test, cmap='Blues')
6
7 plt.grid(False)
8
```



```
In [43]: 1 precision = precision_score(y_test, knn_tuned_preds)
2 recall = recall_score(y_test, knn_tuned_preds)
3
4 print('Precision:', precision)
5 print('Recall:', recall)
```

```
Precision: 0.7813137494900041
Recall: 0.7393822393822393
```

The uniform weighting model performed better on test data.

Reccomendations

- Track seasonal droughts, water conditions and basin levels so we will know what non-functional pumps are un-repairable to better target efforts
- Install more pumps
- Gather more data

Next Steps

- Gather more data
- Re-assess features for significance
- Further tune the model
- Assess more advanced and resource intensive models

In []:

1	
---	--