

1. Overview

In this Python Script, we will apply all the data transformations that were done in the 01_data_preprocessing python script on the test dataset. Moreover, we will also generate predictions of the test dataset with the best logistic regression model that was trained in the 02_model_creation. With all this, we will be able to obtain the predicted values and determine whether a pump will be functional or non-functional

2. Data Understanding

2.1 Data Description

This notebook will use the test dataset given to us from DrivenData called: Test_set_values

2.2 Import Necessary Libraries

```
In [1]: 1 import pickle
        2 import pandas as pd
        3 import re # Import regular expressions library
        4 import warnings
        5 warnings.simplefilter('ignore')
```

2.3 Define global variables

```
In [2]: 1 INPUT_PATH_Test_set_values = "../Data/Test_set_values.csv"
```

2.4 Functions

In [3]:

```
1 def categorize_funder(funder):
2     """
3     Categorizes a funder name into specific groups based on keywords.
4
5     Args:
6     funder (str): A string representing the name of the funder to categori
7
8     Returns:
9     str: A category name representing the type of organization the funder
10
11     This function takes a funder name, converts it to lowercase, removes l
12     and categorizes it into predefined groups like 'Government', 'Religiou
13     'International Aid', 'Private Companies', or 'Individual/Other' based
14     """
15     funder = funder.lower().strip() # convert to lowercase and strip whit
16     if any(x in funder for x in ['government', 'ministry', 'gov', 'minis']):
17         return 'Government'
18     elif any(x in funder for x in ['church', 'muslim', 'mus', 'islamic', 'is
19         return 'Religious Organizations'
20     elif any(x in funder for x in ['ngo', 'foundation', 'fund', 'trust', '
21         return 'NGO'
22     elif any(x in funder for x in ['international', 'internatio', 'un', 'wc
23         return 'International Aid'
24     elif any(x in funder for x in ['ltd', 'company', 'compa', 'group', 'ent
25         return 'Private Companies'
26     else:
27         return 'Individual/Other'
28
```

In [4]:

```
1 def categorize_installer(installer):
2     """
3     Categorizes an installer name into specific groups based on keywords.
4
5     Args:
6     installer (str): A string representing the name of the installer to ca
7
8     Returns:
9     str: A category name representing the type of entity the installer bel
10
11     This function processes an installer name by converting it to lowercas
12     any leading/trailing whitespace. It categorizes the name into predefin
13     'DWE', 'Government', 'Community', 'NGO', 'Private Company', 'Instituti
14     based on specific keywords present in the installer's name. This helps
15     installer data for better analysis and insight extraction.
16     """
17     installer = installer.lower().strip() # convert to lowercase and stri
18     if 'dw' in installer:
19         return 'DWE'
20     elif any(x in installer for x in ['government', 'govt', 'gove']):
21         return 'Government'
22     elif any(x in installer for x in ['resource']):
23         return 'Other'
24     elif any(x in installer for x in ['community', 'villagers', 'village',
25         return 'Community'
26     elif any(x in installer for x in ['ngo', 'unicef', 'foundat']):
27         return 'NGO'
28     elif 'company' in installer or 'contractor' in installer:
29         return 'Private Company'
30     elif any(x in installer for x in ['school', 'schoo', 'church', 'rc']):
31         return 'Institutional'
32     else:
33         return 'Other'
```

```

In [5]: 1 def group_scheme_management(value):
        2     """
        3     Categorizes scheme management types into broader, more generalized groups.
        4
        5     Args:
        6     value (str): A string representing the scheme management type to categorize.
        7
        8     Returns:
        9     str: A generalized category name representing the type of scheme management.
        10
        11     This function takes a specific scheme management type and categorizes it into
        12     more generalized groups such as 'Government', 'Community', 'Private Sector',
        13     'Water Board', or 'Other'. This categorization aids in simplifying the analysis
        14     and understanding of the data by reducing the number of distinct categories,
        15     making trends and patterns more discernible.
        16     """
        17     if value in ['VWC', 'Water authority', 'Parastatal']:
        18         return 'Government'
        19     elif value in ['WUG', 'WUA']:
        20         return 'Community'
        21     elif value in ['Company', 'Private operator']:
        22         return 'Private Sector'
        23     elif value == 'Water Board':
        24         return 'Water Board' # Retain this as a separate category if distinct
        25     else:
        26         return 'Other'

```

```

In [6]: 1 def clean_text(text):
        2     """
        3     Cleans a text string by converting to lowercase, removing non-alphanumeric characters,
        4     and replacing multiple spaces with a single space. If the input is None or NaN, it returns None.
        5
        6     Args:
        7     text (str or NaN): The text to be cleaned; can be a string, numeric, or None/NaN.
        8
        9     Returns:
        10    str or NaN: The cleaned text, with all characters in lowercase, non-alphanumeric characters removed,
        11                and multiple spaces collapsed to a single space, or the original text if it was None/NaN.
        12
        13    This function standardizes a text string by making it lowercase, stripping leading and trailing
        14    whitespace, and then replacing sequences of spaces with a single space, facilitating easier
        15    comparison and analysis. If the input is numeric, it is assumed to be standardized already and is returned
        16    unchanged.
        17    """
        17     if pd.isna(text):
        18         return text
        19     if isinstance(text, (int, float)): # Check if the input is numeric
        20         return text
        21     text = text.lower() # Convert to lowercase
        22     text = ''.join(char for char in text if char.isalpha() or char.isspace())
        23     text = re.sub(r'\s+', ' ', text) # Replace multiple spaces with a single space
        24     return text

```

3. Code

3.1 Import the dataset

```
In [7]: 1 df_predict = pd.read_csv(INPUT_PATH_Test_set_values)
        2 df_predict.head()
```

Out[7]:

| | id | amount_tsh | date_recorded | funder | gps_height | installer | longitude | latitude |
|---|-------|------------|---------------|---------------------------|------------|---------------|-----------|------------|
| 0 | 50785 | 0.0 | 2013-02-04 | Dmdd | 1996 | DMDD | 35.290799 | -4.059696 |
| 1 | 51630 | 0.0 | 2013-02-04 | Government Of Tanzania | 1569 | DWE | 36.656709 | -3.309214 |
| 2 | 17168 | 0.0 | 2013-02-01 | NaN | 1567 | NaN | 34.767863 | -5.004344 |
| 3 | 45559 | 0.0 | 2013-01-22 | Finn Water | 267 | FINN WATER | 38.058046 | -9.418672 |
| 4 | 49871 | 500.0 | 2013-03-27 | Bruder | 1260 | BRUDER | 35.006123 | -10.950412 |

5 rows × 40 columns

```
In [8]: 1 # Modifying columns to keep the structure of the data that the models wher
        2
        3 df_predict['permit'] = df_predict['permit'].replace({True:1})
        4
        5 df_predict['public_meeting'] = df_predict['public_meeting'].replace({True:
```

3.2 Apply the same data transformations on df_predict as the ones done in 00_data_understanding

3.2.1 Applying transformation functions

Column 'funder'

```
In [9]: 1 # Handling NaN values with a filler string like 'Unknown'
2 df_predict['funder'] = df_predict['funder'].fillna('Unknown').astype(str)
3
4 # Apply the mapping function to the 'funder' column
5 df_predict['funder_type'] = df_predict['funder'].apply(categorize_funder)
6
7 # Check the categorized data
8 print(df_predict['funder_type'].value_counts())
```

```
Individual/Other          9955
Government                2438
International Aid        2093
Religious Organizations   329
NGO                      29
Private Companies          6
Name: funder_type, dtype: int64
```

Column 'installer'

```
In [10]: 1 # Handling NaN values with a filler string like 'Unknown'
2 df_predict['installer'] = df_predict['installer'].fillna('Unknown').astype(str)
3
4 # Apply the mapping function to the 'installer' column
5 df_predict['installer_type'] = df_predict['installer'].apply(categorize_installer)
6
7 # Now you can check your categorized data
8 print(df_predict['installer_type'].value_counts())
```

```
Other          8480
DWE            4537
Government     926
Community      599
Institutional  185
NGO            93
Private Company 30
Name: installer_type, dtype: int64
```

Column 'scheme_management_grouped'

```
In [11]: 1 # Apply the grouping function to the 'scheme_management' column
2 df_predict['scheme_management_grouped'] = df_predict['scheme_management'].apply(grouping_function)
3
4 # Check the new value counts to see the grouped data
5 print(df_predict['scheme_management_grouped'].value_counts(normalize=True))
```

```
Government    0.699663
Community     0.131852
Other         0.083838
Water Board   0.048081
Private Sector 0.036566
Name: scheme_management_grouped, dtype: float64
```

3.2.2 Converting data types

```
In [12]: 1 # Converting 'construction_year' to object
        2 df_predict['construction_year'] = df_predict['construction_year'].astype('object')
```

```
In [13]: 1 df_predict.columns
```

```
Out[13]: Index(['id', 'amount_tsh', 'date_recorded', 'funder', 'gps_height',
               'installer', 'longitude', 'latitude', 'wpt_name', 'num_private',
               'basin', 'subvillage', 'region', 'region_code', 'district_code', 'lg
a',
               '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', 'funder_type', 'installer_type',
               'scheme_management_grouped'],
              dtype='object')
```

3.2.3 Drop unnecessary columns

```
In [14]: 1 drop_column_list = ['scheme_name', 'num_private', 'wpt_name', 'subvillage',
        2                       'extraction_type', 'management', 'payment', 'water_quality',
        3                       'waterpoint_type_group', 'date_recorded', 'funder', 'installer',
        4                       'longitude', 'latitude', 'region_code', 'district_code', 'lg_a']
```

```
In [15]: 1 df_predict = df_predict.drop(drop_column_list, axis=1)
```

3.2.3 Cleaning the data set

```
In [16]: 1 # Apply the cleaning function to each object-type column in the DataFrame
        2 for col in df_predict.select_dtypes(include='object').columns:
        3     df_predict[col] = df_predict[col].apply(clean_text)
```

3.2 Fillna with the modes calculated in 01_data_preprocessing

```
In [17]: 1 (df_predict.isna().sum()/len(df_predict))*100
```

```
Out[17]: id                0.000000
amount_tsh              0.000000
gps_height              0.000000
basin                   0.000000
region                  0.000000
population              0.000000
public_meeting          5.528620
permit                  4.962963
extraction_type_class   0.000000
management_group        0.000000
payment_type            0.000000
quality_group           0.000000
quantity_group          0.000000
source_type             0.000000
waterpoint_type         0.000000
funder_type             0.000000
installer_type          0.000000
scheme_management_grouped 0.000000
dtype: float64
```

From the python script 01_data_preprocessing we know that public_meeting_mode is 1.0 and the permit_mode is 1.0. So we are going to directly fill the NaNs of public_meeting and of permit with the value 1.0

Fillna in column 'public_meeting'

```
In [18]: 1 df_predict['public_meeting'].fillna(1.0, inplace=True)
```

Fillna in column 'permit'

```
In [19]: 1 df_predict['permit'].fillna(1.0, inplace=True)
```

Let's check that there are no more null-values left


```
In [20]: 1 (df_predict.isna().sum()/len(df_predict))*100
```

```
Out[20]: id                0.0
amount_tsh              0.0
gps_height              0.0
basin                  0.0
region                 0.0
population             0.0
public_meeting         0.0
permit                0.0
extraction_type_class  0.0
management_group       0.0
payment_type           0.0
quality_group          0.0
quantity_group         0.0
source_type            0.0
waterpoint_type        0.0
funder_type            0.0
installer_type         0.0
scheme_management_grouped 0.0
dtype: float64
```

3.3 Doing target encoder on the categorical columns

Let's apply a one hot encoder for the categorical columns that have 6 or less categories

```
In [21]: 1 # Capture categorical columns from X_train for encoding
2 categorical_columns = df_predict.select_dtypes(include=['object', 'categor
3
4
5 # Encoding the categorical columns in df_predict
6 for col in categorical_columns:
7     if df_predict[col].nunique() <= 6:
8         # Apply OneHotEncoder for columns with 6 or fewer unique values
9         df_predict = pd.get_dummies(df_predict, columns=[col], drop_first=
```

Let's call in the saved fits (for the categorical columns that have more than 6 categories) applied to the categorical columns in the 01_data_preprocessing script

In [22]: 1 df_predict.columns

```
Out[22]: Index(['id', 'amount_tsh', 'gps_height', 'basin', 'region', 'population',
               'extraction_type_class', 'payment_type', 'source_type',
               'waterpoint_type', 'installer_type', 'public_meeting_1', 'permit_1',
               'management_group_other', 'management_group_parastatal',
               'management_group_unknown', 'management_group_usergroup',
               'quality_group_fluoride', 'quality_group_good', 'quality_group_milky',
               'quality_group_salty', 'quality_group_unknown', 'quantity_group_enoug
h',
               'quantity_group_insufficient', 'quantity_group_seasonal',
               'quantity_group_unknown', 'funder_type_individualother',
               'funder_type_international aid', 'funder_type_ngo',
               'funder_type_private companies', 'funder_type_religious organization
s',
               'scheme_management_grouped_government',
               'scheme_management_grouped_other',
               'scheme_management_grouped_private sector',
               'scheme_management_grouped_water board'],
              dtype='object')
```

```
In [23]: 1 # Column 'basin'
2 basin_pickle = pickle.load(open('model_objects/basin_target_encoder.pkl',
3 df_predict['basin'] = basin_pickle.transform(df_predict['basin'])
4
5 # Column 'extraction_type_class'
6 extraction_type_class_pickle = pickle.load(open('model_objects/extraction_
7 df_predict['extraction_type_class'] = extraction_type_class_pickle.transf
8
9 # Column 'installer_type'
10 installer_type_pickle = pickle.load(open('model_objects/installer_type_tar
11 df_predict['installer_type'] = installer_type_pickle.transform(df_predict[
12
13 # Column 'payment_type'
14 payment_type_pickle = pickle.load(open('model_objects/payment_type_target_
15 df_predict['payment_type'] = payment_type_pickle.transform(df_predict['pay
16
17 # Column 'region_target'
18 region_target_pickle = pickle.load(open('model_objects/region_target_encoc
19 df_predict['region'] = region_target_pickle.transform(df_predict['region'])
20
21 # Column 'source_type'
22 source_type_pickle = pickle.load(open('model_objects/source_type_target_er
23 df_predict['source_type'] = source_type_pickle.transform(df_predict['sourc
24
25 # Column 'waterpoint_type'
26 waterpoint_type_pickle = pickle.load(open('model_objects/waterpoint_type_t
27 df_predict['waterpoint_type'] = waterpoint_type_pickle.transform(df_predic
```

3.4 Dealing with numerical columns

Let's call in the saved fits applied to the numerical columns in the 01_data_preprocessing script

In [24]: 1 df_predict.columns

Out[24]: Index(['id', 'amount_tsh', 'gps_height', 'basin', 'region', 'population',
'extraction_type_class', 'payment_type', 'source_type',
'waterpoint_type', 'installer_type', 'public_meeting_1', 'permit_1',
'management_group_other', 'management_group_parastatal',
'management_group_unknown', 'management_group_usergroup',
'quality_group_fluoride', 'quality_group_good', 'quality_group_milky',
'quality_group_salty', 'quality_group_unknown', 'quantity_group_enough',
'quantity_group_insufficient', 'quantity_group_seasonal',
'quantity_group_unknown', 'funder_type_individualother',
'funder_type_international aid', 'funder_type_ngo',
'funder_type_private companies', 'funder_type_religious organization',
'scheme_management_grouped_government',
'scheme_management_grouped_other',
'scheme_management_grouped_private sector',
'scheme_management_grouped_water board'],
dtype='object')

```
In [25]: 1 # Capture numerical columns
2 numerical_columns = df_predict.select_dtypes(include=['int64', 'float64'])
3
4 # Let's also drop column 'id' from the numerical_columns as they don't serve
5 numerical_columns = numerical_columns.drop('id')
6
7 # The desired order of columns according to when it was fitted in notebook
8 desired_order = [
9     'amount_tsh', 'gps_height', 'population', 'basin', 'region',
10    'extraction_type_class', 'payment_type', 'source_type',
11    'waterpoint_type', 'installer_type'
12 ]
13
14 # Numerical Columns
15 numerical_columns_pickle = pickle.load(open('model_objects/numerical_columns.pkl', 'rb'))
16 df_predict[desired_order] = numerical_columns_pickle.transform(df_predict[desired_order])
```

```
In [26]: 1 df_predict.columns
```

```
Out[26]: Index(['id', 'amount_tsh', 'gps_height', 'basin', 'region', 'population',  
              'extraction_type_class', 'payment_type', 'source_type',  
              'waterpoint_type', 'installer_type', 'public_meeting_1', 'permit_1',  
              'management_group_other', 'management_group_parastatal',  
              'management_group_unknown', 'management_group_usergroup',  
              'quality_group_fluoride', 'quality_group_good', 'quality_group_milky',  
              'quality_group_salty', 'quality_group_unknown', 'quantity_group_enoug  
h',  
              'quantity_group_insufficient', 'quantity_group_seasonal',  
              'quantity_group_unknown', 'funder_type_individualother',  
              'funder_type_international aid', 'funder_type_ngo',  
              'funder_type_private companies', 'funder_type_religious organization  
s',  
              'scheme_management_grouped_government',  
              'scheme_management_grouped_other',  
              'scheme_management_grouped_private sector',  
              'scheme_management_grouped_water board'],  
              dtype='object')
```

3.5 Apply the Decision Tree Classifier created in 02_model_creation

```
In [27]: 1 df_predict
```

```
Out[27]:
```

| | id | amount_tsh | gps_height | basin | region | population | extraction_type_class | p |
|-------|-------|------------|------------|-----------|-----------|------------|-----------------------|---|
| 0 | 50785 | -0.100621 | 1.915327 | -0.540016 | -0.633090 | 0.299241 | 2.617222 | |
| 1 | 51630 | -0.100621 | 1.299135 | -0.569630 | -1.180350 | 0.254822 | -0.521411 | |
| 2 | 17168 | -0.100621 | 1.296248 | -0.540016 | 0.663779 | 0.677863 | 2.617222 | |
| 3 | 45559 | -0.100621 | -0.579749 | 2.497921 | 2.407560 | 0.149062 | 2.617222 | |
| 4 | 49871 | 0.055600 | 0.853225 | 2.497921 | -0.015500 | -0.252826 | -0.521411 | |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 14845 | 39307 | -0.100621 | -0.915985 | 0.697368 | 0.103890 | -0.337435 | 1.165688 | |
| 14846 | 18990 | 0.211821 | -0.965049 | -0.569630 | 0.234762 | 5.881260 | -0.463637 | |
| 14847 | 28749 | -0.100621 | 1.164929 | -0.540016 | 0.663779 | 0.043302 | -0.521411 | |
| 14848 | 33492 | -0.100621 | 0.475139 | -1.230325 | -0.015500 | -0.062458 | -0.521411 | |
| 14849 | 68707 | -0.100621 | -0.270931 | -1.230325 | -0.015500 | -0.295131 | -0.521411 | |

14850 rows × 35 columns

```
In [28]: 1 # Loading the pickle for the best Decision Tree Classifier
2 best_tree_pickle = pickle.load(open('model_objects/best_tree.pkl', 'rb'))
3
4 # Let's drop the 'id' column
5 df_predict_copy = df_predict.drop('id', axis=1)
```

```
In [29]: 1 df_predict_copy.rename(columns={'public_meeting_1': 'public_meeting_1.0', '
2
3 # Let's reorder the columns
4 df_predict_copy = df_predict_copy[list(best_tree_pickle.feature_names_in_)
```

```
In [30]: 1 # Decision Tree Classifier
2 df_predict['status_group'] = best_tree_pickle.predict_proba(df_predict_cop
```

```
In [31]: 1 # Apply a threshold to the probabilities of status_group to determine to w
2 df_predict['status_group_class'] = df_predict['status_group'].map(lambda x
```

```
In [32]: 1 df_predict[['id', 'status_group', 'status_group_class']]
```

Out[32]:

| | id | status_group | status_group_class | |
|--|-------|--------------|--------------------|----------------|
| | 0 | 50785 | 0.348329 | Functional |
| | 1 | 51630 | 0.000000 | Functional |
| | 2 | 17168 | 0.000000 | Functional |
| | 3 | 45559 | 0.989969 | Non-functional |
| | 4 | 49871 | 0.179576 | Functional |
| | ... | ... | ... | ... |
| | 14845 | 39307 | 0.866786 | Non-functional |
| | 14846 | 18990 | 0.255741 | Functional |
| | 14847 | 28749 | 0.038544 | Functional |
| | 14848 | 33492 | 0.698860 | Non-functional |
| | 14849 | 68707 | 0.994873 | Non-functional |

14850 rows × 3 columns

4. Export the data

```
In [33]: 1 df_predict[['id', 'status_group_class']].to_excel('Final_results.xlsx', ir
```