

Section-V- Self Case Study-1_Pump it up-Data Mining the Water Table

7. Creating a function for predicting the result of test dataset

```
In [16]: import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import re
import time
import warnings
warnings.filterwarnings("ignore")
import numpy as np
from nltk.corpus import stopwords
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
import sys
from category_encoders import OneHotEncoder
from sklearn.preprocessing import MinMaxScaler
from category_encoders import TargetEncoder, LeaveOneOutEncoder, WOEEncoder
from sklearn.preprocessing import RobustScaler
from sklearn.metrics import accuracy_score, balanced_accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from xgboost import XGBClassifier
from prettytable import PrettyTable
import category_encoders as ce
from imblearn.over_sampling import SMOTE
from xgboost import XGBClassifier
from sklearn.metrics import f1_score
from sklearn.metrics import roc_auc_score
import joblib
```

```
In [17]: df_test = pd.read_csv("test.csv")
```

```

In [19]: def DMWT_preprocessor(df):
    df['funder'].fillna(value='Undefined',inplace=True)
    df['funder'].replace(to_replace = '0', value = 'Undefined' , inplace=True) #replacing '0' & missing values with 'Undefined'
    top_30_funders = ['Government Of Tanzania','Undefined','Danida','Hesawa','Rwssp','World Bank','Kkkt','World Vision',
        'Unicef','Tasaf','District Council','Dhv','Private Individual','Dwsp','Norad','Germany Republi',
        'Tcrs','Ministry Of Water','Water','Dwe','Netherlands','Hifab','Adb','Lga','Amref','Fini Water',
        'Oxfam','Wateraid','Rc Church','Isf']
    df.loc[~df["funder"].isin(top_30_funders), "funder"] = "other"

    df['installer'].fillna(value='Undefined',inplace=True)
    df['installer'].replace(to_replace = '0', value = 'Undefined' , inplace=True) #replacing '0' category with 'Undefined'

    df['installer'].replace(to_replace = ("Gove","Gover","GOVERN", "GOVERN", "GOVERNME", "Governmen","Government",
        "GOVER"), value = "Government", inplace=True)
    df['installer'].replace(to_replace = ("RW","RWE","RWE /Community","RWE Community","RWE/ Community","RWE/Community",
        "RWE/DWE","RWE/TCRS","RWEDWE","RWET/WESA"), value = "RWE", inplace=True)
    df['installer'].replace(to_replace = ("Commu", "Communit", "Community", "COMMUNITY BANK",
        "Comunity"), value = "Community", inplace=True)
    df['installer'].replace(to_replace = ("Danda", "DANIAD", "Danid", "DANIDA", "DANIDA CO", "DANIDS", "DANNIDA",
        "DANID"), value = "DANIDA", inplace=True)
    df['installer'].replace(to_replace = ("Cebtral Government", "Cental Government", "Centr", "Centra Government", "Centra govt",
        "Central government", "Central govt", "Cetral government /RC", "Tanzania Government",
        "TANZANIAN GOVERNMENT"), value = "Central Government", inplace=True)
    df['installer'].replace(to_replace = ("COUN", "Counc", "Council", "Distri", "District Council",
        "District Community j", "District Counci", "District council",
        "District Council"), value = "District Council", inplace=True)
    df['installer'].replace(to_replace = ("Hesawa", "HESAW", "Hesewa", "HESAWA"), value = 'HESAWA' , inplace=True)
    df['installer'].replace(to_replace = ("World Division", "World Visiin", "World vision", "World Vission",
        "World Vision"), value = 'World Vision' , inplace=True)
    df['installer'].replace(to_replace = ("Distric Water Department", "District Water Department",
        "District water depar", "District water department",
        "Water Department"), value = 'District water department' , inplace=True)
    df['installer'].replace(to_replace = ("FINN WATER", "FinW", "FinWate", "FinWater", "Fini water",
        "Fini Water" ), value = 'Fini Water' , inplace=True)
    df['installer'].replace(to_replace = ("RC", "RC .Church", "RC C", "RC Ch", "RC Churc", "RC Church",
        "RC CHURCH BROTHER", "RC church/CEFA", "RC church/Central Gover",
        "RCchurch/CEFA", "RC CHURCH"), value = "RC Church" , inplace=True)
    df['installer'].replace(to_replace = ("Villa", "VILLAGER", "Villagerd", "Villagers", "Villages", "Village Council",
        "Villi", "villigers"), value = "Villagers" , inplace=True)

    top_30_installer = ["DWE", "Undefined", "Government", "DANIDA", "Community", "HESAWA", "RWE", "District Council",
        "Central Government", "KKKT", "TCRS", "World Vision", "CES", "Fini Water", "RC Church", "LGA",
        "WEDECO", "TASAF", "AMREF", "TWESA", "WU", "Dmdd", "ACRA", "Villagers", "SEMA", "DW", "OXFAM", "Da",
        "Idara ya maji", "UNICEF"]

```

```

df.loc[~df["installer"].isin(top_30_installer), "installer"] = "other"

df['longitude'].replace(to_replace = 0 , value = 34.07742669202832 , inplace=True)
df['population'].replace(to_replace = 0, value = 281.087167 , inplace=True)
df['construction_year'].replace(to_replace = 0, value = 1996, inplace=True)
df['date_recorded'] = pd.to_datetime(df['date_recorded']) #converting dates to 'datetime' datatype
df['operational_year'] = df.date_recorded.dt.year - df.construction_year
df.operational_year.head(5)
df.loc[df['operational_year']<0, 'operational_year'] = 0

df.drop(columns=["construction_year", "date_recorded", "extraction_type", "extraction_type_class", 'payment_type',
                "quality_group", "quantity_group", "source_type", "source_class", "waterpoint_type_group", 'permit',
                "scheme_management", 'id', 'amount_tsh', 'ward', 'wpt_name', 'num_private', 'subvillage', 'region_code',
                'public_meeting', 'recorded_by', 'management_group', 'scheme_name'], inplace=True)

#encoding target variables manually

numerical_features = ['gps_height', 'longitude', 'latitude', 'district_code', 'population', 'operational_year']
categorical_features = ['funder', 'installer', 'basin', 'region', 'lga', 'extraction_type_group', 'management',
                        'payment', 'water_quality', 'quantity', 'source', 'waterpoint_type']

scaler = joblib.load('scaler.pkl')
df[numerical_features] = scaler.transform(df[numerical_features])

encoder = joblib.load('encoder.pkl')
df[categorical_features] = encoder.transform(df[categorical_features])

return df

```

```

In [20]: def DMWT_prediction(df):
df = DMWT_preprocessor(df)
model = joblib.load('model.pkl')

y_pred = model.predict(df)
y_pred_train_proba = model.predict_proba(df)

return y_pred, y_pred_train_proba

```

```
In [21]: DMWT_prediction(df_test)
```

```
Out[21]: (array([1, 1, 1, ..., 1, 1, 0], dtype=int64),  
          array([[0.24287999, 0.72248787, 0.03463215],  
                [0.2529064 , 0.5026612 , 0.2444323 ],  
                [0.03922484, 0.878787 , 0.08198813],  
                ...,  
                [0.22526748, 0.7267688 , 0.04796375],  
                [0.06310199, 0.8525588 , 0.08433919],  
                [0.9881142 , 0.01067844, 0.00120742]], dtype=float32))
```

```
In [ ]:
```