PROJECT DEVELOPMENT PHASE - SPRINT III

Assignment Date	09-11-2022
Team ID	PNT2022TMID06599
Project Name	Efficient Water Quality Analysis and Prediction using Machine Learning
Maximum Marks	8 Mark

Train and Develop the Model

Data Collection:

import numpy as np import pandas as pd import seaborn as sns import matplotlib.pyplot as plt import matplotlib as mpl import matplotlib.patches as patches from matplotlib.patches import ConnectionPatch from collections import OrderedDict from matplotlib.gridspec import GridSpec %matplotlib inline

```
df = pd.read_csv('Final.csv')
df
```

Exploratory Data Analysis:

```
df.shape

df.isnull().sum()

df.info()

df.describe()

df.fillna(df.mean(), inplace=True)
df.isnull().sum()
```

```
df.Potability.value counts()
sns.countplot(df['Potability'])
plt.show()
sns.distplot(df['ph'])
plt.show()
df.hist(figsize=(14,14))
plt.show()
plt.figure(figsize=(13,8))
sns.heatmap(df.corr(),annot=True,cmap='terrain')
plt.show()
df.boxplot(figsize=(14,7))
X = df.drop('Potability',axis=1)
Y= df['Potability']
from sklearn.model selection import train test split
X train, X test, Y train, Y test = train test split(X,Y, test size=0.2,
random state=101,shuffle=True)
Train Decision Tree Classifier and check accuracy:
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy score, confusion matrix, classification report
dt=DecisionTreeClassifier(criterion='gini', min samples split= 10, splitter='best')
dt.fit(X train,Y train)
prediction=dt.predict(X test)
print(f"Accuracy Score = {accuracy score(Y test,prediction)*100}")
print(f''Confusion Matrix = \n \{confusion matrix(Y test, prediction)\}'')
print(f"Classification Report = \n {classification report(Y test, prediction)}")
```

```
res = dt.predict([[7.408985467,0.57139761,40,6.505923139,311.4526625,504.1459941, 11.53214401,81.10693773,3.772420928,0.0,100,0.0,16.5,0.0,11.24]])[0] res
```

Apply Hyper Parameter Tuning:

```
from sklearn.model selection import RepeatedStratifiedKFold
from sklearn.model selection import GridSearchCV
# define models and parameters
model = DecisionTreeClassifier()
criterion = ["gini", "entropy"]
splitter = ["best", "random"]
min samples split = [2,4,6,8,10,12,14]
# define grid search
grid = dict(splitter=splitter, criterion=criterion,
min samples split=min samples split)
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
grid search dt = GridSearchCV(estimator=model, param grid=grid, n jobs=-1,
cv=cv,
                scoring='accuracy',error score=0)
grid search dt.fit(X train, Y train)
print(f''Best: {grid search dt.best score :.3f} using
{grid search dt.best params }")
means = grid search dt.cv results ['mean test score']
stds = grid search dt.cv results ['std test score']
params = grid search dt.cv results ['params']
```

```
for mean, stdev, param in zip(means, stds, params):
    print(f"{mean:.3f} ({stdev:.3f}) with: {param}")

print("Training Score:",grid_search_dt.score(X_train, Y_train)*100)
print("Testing Score:", grid_search_dt.score(X_test, Y_test)*100)
```

Modelling:

```
df.head(20)
df.tail(5)
df['Potability'].value_counts().to_frame()
df_filtered = df[df['Turbidity'].isin(["1,2,3,4,5,6,7,8,9,10"])]
print(df_filtered.head(15))
print(df_filtered.shape)
```

Model Evaluation

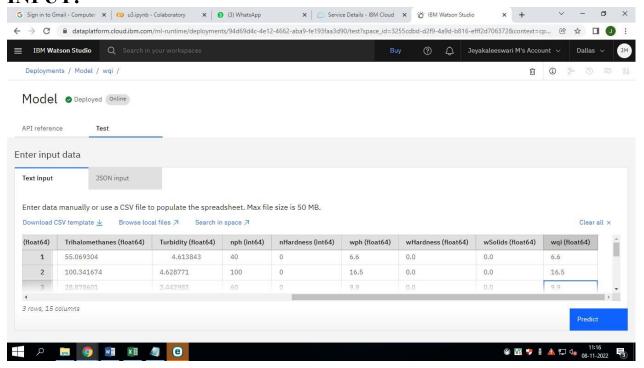
```
from sklearn.metrics import r2_score from sklearn.metrics import mean_absolute_error from sklearn.metrics import mean_squared_error print('R Squared=',r2_score(X_train,Y_test)) print('MAE=',mean_absolute_error(X_train,Y_test)) print('MSE=',mean_squared_error(X_train,Y_test)) print(joblib joblib.dump(dt, 'classifier.pkl')
```

!pip install -U ibm-watson-machine-learning

```
from ibm watson machine learning import APIClient
import json
import numpy as np
wml credentials =
{"apikey":"nFFWACn7pVNTQWlnb7pusoXVa63g0vFEq 8Y2x2pxZSE","url":
          "https://us-south.ml.cloud.ibm.com" }
wml client = APIClient(wml credentials)
wml client.spaces.list()
SPACE ID = "3255cdbd-d2f9-4a9d-b816-efff2d706372"
wml client.set.default space(SPACE ID)
wml client.software specifications.list(500)
import sklearn
sklearn. version
MODEL NAME = 'wqi'
DEPLOYMENT NAME = 'Model'
DEMO MODEL = dt
# Set Python Version
software spec uid =
wml client.software specifications.get id by name('runtime-22.1-py3.9')
# Setup model meta
model props = {
  wml_client.repository.ModelMetaNames.NAME: MODEL_NAME,
  wml_client.repository.ModelMetaNames.TYPE: 'scikit-learn_1.0',
  wml_client.repository.ModelMetaNames.SOFTWARE_SPEC_UID:
software spec uid
```

```
SAVE THE MODEL:
#Save model
model details =
  wml client.repository.store model(model=DEM
  O MODEL, meta props=model props,
  training data=X train,
  training target=Y train
model details
model_id = wml_client.repository.get_model_id(model_details)
model id
# Set meta
deployment props = {
wml client.deployments.ConfigurationMetaNames.NAME:DEPLOYMENT NA
ME,
  wml client.deployments.ConfigurationMetaNames.ONLINE: {}
DEPLOY:
# Deploy
deployment =
  wml client.deployments.create(artifact uid
  =model id, meta props=deployment props
)
```

INPUT:



OUTPUT:

