# Analyzing Water Quality Data using Distributed Cloud Services and Spark

Abdoulaye Diallo
*New York Institute Of Technology*
New York, NY.
adiall04@nyit.edu

## 1. INTRODUCTION:

Having clean water is crucial for our health. Determining if water is safe involves many factors. Since water quality depends on the amounts of certain elements in it, we can create a model to check water quality based on these elements.

## 2. PACKAGES

```python
1  # Importing necessary libraries
2  import pyspark
3  from pyspark.sql import SparkSession
4  from pyspark.sql.functions import count, when, col
5  from pyspark.sql.functions import col
6  from pyspark.ml.feature import VectorAssembler
7
8  import numpy as np
9  import pandas as pd
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 import plotly.express as px
13 import missingno as msno
```

### PREPARING THE DATA FOR MODELLING

```python
7]:  1  # Import statement for splitting data into training and testing sets
     2  from sklearn.model_selection import train_test_split
     3
     4  # Import statements for various scoring metrics
     5  from sklearn.metrics import precision_score, accuracy_score, mean_absolute_error, mean_squared_error
     6
     7  # Import statements for different classifier models
     8  from sklearn.linear_model import LogisticRegression, RidgeClassifier, SGDClassifier, Perceptron, PassiveAggressiveC
     9  from sklearn.svm import SVC, NuSVC
     10 from sklearn.tree import DecisionTreeClassifier
     11 from sklearn.naive_bayes import GaussianNB, BernoulliNB
     12 from sklearn.neighbors import NearestCentroid, KNeighborsClassifier
     13 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
     14
     15 # Import statement for XGBoost classifier
     16 from xgboost import XGBClassifier
     17
     18 #Hyper-Parameter-Tuning
     19 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
     20
```

```python
n [351]:  1  # Import necessary functions from sklearn.metrics for model evaluation
          2  from sklearn.metrics import confusion_matrix          # For calculating confusion matrix
          3  from sklearn.metrics import classification_report     # For generating classification report
          4  from sklearn.metrics import precision_score, recall_score, f1_score  # For calculating precision, recall, and F1-sc
          5  from sklearn.metrics import roc_curve, auc
          6
```

## 3. BACKGROUND AND MOTIVATION:

Having access to safe drinking water is not only crucial for maintaining good health, but it's also recognized as a fundamental human right and a key element of well-rounded health protection policies. This significance extends across national, regional, and local contexts, serving as a vital concern for both health and developmental matters. Demonstrated in some areas, investments made in water supply and sanitation have proven to bring about a positive net economic

outcome. This is due to the fact that the economic gains resulting from reduced health issues and lowered healthcare expenses outweigh the costs associated with implementing these interventions.

4. **DATA LAKE VS DATA STORE:**

**Definitions:**

**Data warehouse:**

A data warehouse is a specialized system for managing data, strategically crafted to facilitate and empower business intelligence (BI) projects, particularly analytics. These repositories are purpose-built for executing queries and conducting analyses, often housing extensive historical data. The content of a data warehouse typically originates from diverse origins, encompassing elements like application log files and transactional applications.

**Data lake:**

A data lake serves as a centralized storage hub where you can house all your structured and unstructured data, regardless of its volume. It permits the storage of data in its raw state, eliminating the need for upfront structuring. This setup facilitates a spectrum of analytical activities, spanning from creating dashboards and visualizations to engaging in substantial data processing, real-time analytics, and leveraging machine learning for informed decision-making. I hosted the data in **AWS S3,** made it public ang granted a user full S3 Access which permitted me to get a **secret** and **access keys** that I used in a jupyter notebook to get the data using spark. However, choosing between a Data Lake and a Data Warehouse depends on various factors such as the company goals and requirements, the nature of the data itself, the goals of the analysis and many factors. Let's consider some specific aspects to our given dataset. In this case, Lets lean more towards a data warehouse, even though S3 is mostly used as storage due to its flexibility, scalability and distributed across three **Availability zones.**

- ❖ **Data Variety and Flexibility:** If the water quality dataset includes a wide range of data types, such as sensor readings, images, text reports, and more, a Data Lake might be a better choice. Data Lakes excel at handling diverse and unstructured data types without the need for extensive upfront transformation.
- ❖ **Data Volume:** If the water quality dataset contains massive volumes of data, particularly if it's generated from multiple sources like sensors and monitoring stations, a Data Lake might be more scalable. Data Lakes can handle the raw and unaggregated data efficiently.
- ❖ **Data Exploration and Research:** If the primary goal is to explore the water quality data, conduct research, and gain insights into patterns and anomalies, a Data Lake can provide more flexibility. Researchers can directly access the raw data and perform exploratory analysis using various tools.
- ❖ **Regulatory Reporting and Compliance:** If the focus is on structured reporting, compliance, and regular business intelligence reporting related to water quality metrics, a Data Warehouse might be more suitable. Data Warehouses allow you to structure and model the data in a way that aligns with reporting requirements.

- ❖ **Data Quality and Governance:** If ensuring data quality and governance is crucial, a Data Warehouse can provide better control over data transformation and consistency. A well-structured Data Warehouse ensures that only high-quality, relevant data is used for analysis.
- ❖ **Future-Proofing and Agility:** If the water quality dataset is expected to evolve with new data sources, sensors, and changing requirements, a Data Lake offers future-proofing by storing raw data that can be transformed and structured as needs change.
- ❖ **Hybrid Approach:** Depending on the use cases, a hybrid approach could also be considered. Raw water quality data can be stored in a Data Lake for exploratory analysis, and refined and structured data can be moved to a Data Warehouse for reporting and structured analytics.

5. **DATA CONTENT:**
   - ❖ **ph**: It is an important parameter in evaluating the acid–base balance of water. It is also the indicator of acidic or alkaline condition of water status.
   - ❖ **Hardness**: Mainly caused by calcium and magnesium salts. These salts are dissolved from geologic deposits through which water travels. The length of time water is in contact with hardness producing material helps determine how much hardness there is in raw water. Hardness was originally defined as the capacity of water to precipitate soap caused by Calcium and Magnesium.
   - ❖ **Solids**: Water has the ability to dissolve a wide range of inorganic and some organic minerals or salts such as potassium, calcium, sodium, bicarbonates, chlorides, magnesium, sulfates etc. These minerals produced an unwanted taste and diluted color in the appearance of water. This is the important parameter for the use of water.
   - ❖ **Chlorine**: Chlorine and chloramine are the major disinfectants used in public water systems. Chloramines are most commonly formed when ammonia is added to chlorine to treat drinking water.
   - ❖ **Sulfate**: Sulfates are naturally occurring substances that are found in minerals, soil, and rocks. They are present in ambient air, groundwater, plants, and food. The principal commercial use of sulfate is in the chemical industry.
   - ❖ **Conductivity**: Pure water is not a good conductor of electric current rather it's a good insulator. Increase in ions concentration enhances the electrical conductivity of water. Generally, the amount of dissolved solids in water determines the electrical conductivity
   - ❖ **Organic_carbon**: Total Organic Carbon (TOC) in source waters comes from decaying natural organic matter (NOM) as well as synthetic sources. TOC is a measure of the total amount of carbon in organic compounds in pure water.
   - ❖ **Trihalomethanes**: THMs are chemicals which may be found in water treated with chlorine. The concentration of THMs in drinking water varies according to the level of organic material in the water, the amount of chlorine required to treat the water, and the temperature of the water that is being treated.

❖ **Turbidity**: The turbidity of water depends on the quantity of solid matter present in the suspended state. It is a measure of light emitting properties of water and the test is used to indicate the quality of waste discharge with respect to colloidal matter.

❖ **Potability**: Indicates if water is safe for human consumption where 1 means Potable and 0 means Not potable.

## 6. EXPLORATORY DATA ANALYSIS:

A. There are 3276 Observations.

```
There are total 3276 row, Let print first 2 data rows:
Data overview
root
 |-- ph: double (nullable = true)
 |-- Hardness: double (nullable = true)
 |-- Solids: double (nullable = true)
 |-- Chloramines: double (nullable = true)
 |-- Sulfate: double (nullable = true)
 |-- Conductivity: double (nullable = true)
 |-- Organic_carbon: double (nullable = true)
 |-- Trihalomethanes: double (nullable = true)
 |-- Turbidity: double (nullable = true)
 |-- Potability: integer (nullable = true)

Columns overview
```

ıt[307]:

| | Column Name | Data type |
|---|---|---|
| 0 | ph | double |
| 1 | Hardness | double |
| 2 | Solids | double |
| 3 | Chloramines | double |
| 4 | Sulfate | double |
| 5 | Conductivity | double |
| 6 | Organic_carbon | double |
| 7 | Trihalomethanes | double |
| 8 | Turbidity | double |
| 9 | Potability | int |

B. The missing values are mostly on the sulfate (781), 'ph' (491) and 'Trihalomethanes(162) columns.

*Apply necessary data cleaning steps using Spark's DataFrame transformations..*

n [308]:
```
1 # Count missing values in each column
2 missing_counts = df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns])
3 missing_counts.show()
4 plt.show()
```

| ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|
| 491 | 0 | 0 | 0 | 781 | 0 | 0 | 162 | 0 | 0 |

C. To solve the missing values, first create a list of the numeric columns (ph, Sulfate, and Trihalomethanes)  where we have the missing values. Then , loop over that list and for each column, calculate the mean value and subsequently, replace missing values with the mean value. This will help us prepare the data for analysis.

D. Now we check if there is any missing value. Effectively, we find that there is no missing value.
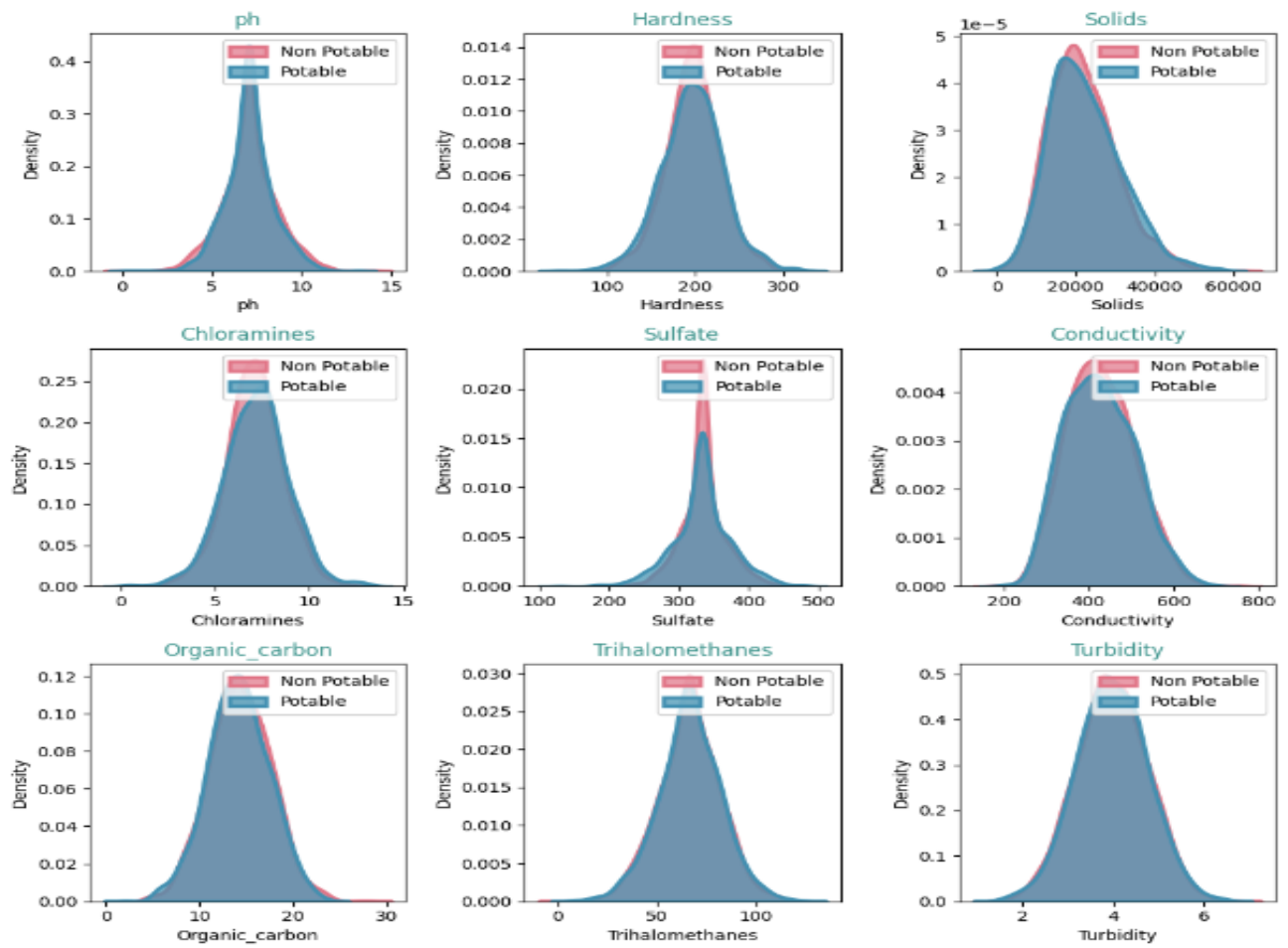
```
In [310]:   1  # check if there are any missing values
            2  missing_counts = df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns])
            3  missing_counts.show()
            4  #Get the number of rows
            5  print(f'There are total {df.count()} row, Let print first 2 data rows:')
```

```
+---+--------+------+----------+-------+------------+--------------+---------------+---------+----------+
| ph|Hardness|Solids|Chloramines|Sulfate|Conductivity|Organic_carbon|Trihalomethanes|Turbidity|Potability|
+---+--------+------+----------+-------+------------+--------------+---------------+---------+----------+
|  0|       0|     0|         0|      0|           0|             0|              0|        0|         0|
+---+--------+------+----------+-------+------------+--------------+---------------+---------+----------+

There are total 3276 row, Let print first 2 data rows:
```
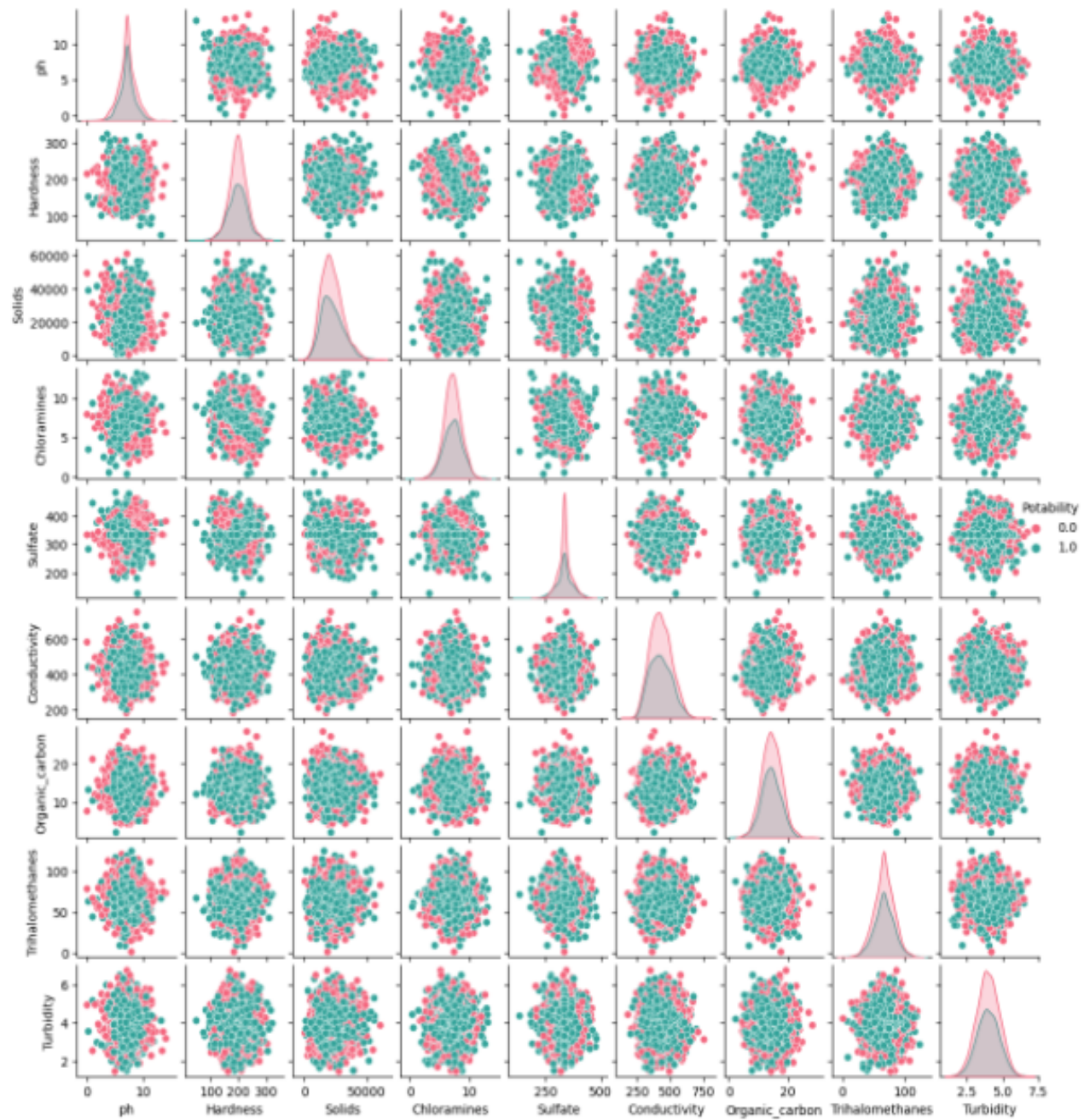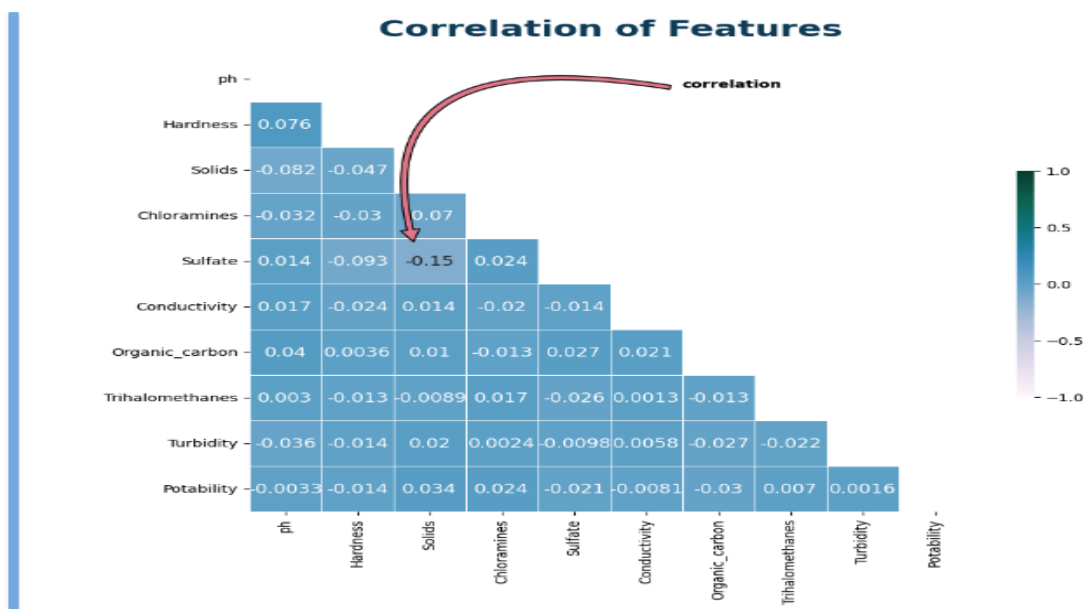
7. **DATA VISUALIZATION:** Since, there are no missing values, let's visualize it to get a more intuitive way to understand the underlying patterns, trends, and relationships of our data.

# Distribution of Features

Water Quality

**Correlation of Features**

## 8. MODELIZATION:

I will  80% for training and 20% for testing the data

```
329]:   1  X = pandas_df.drop("Potability", axis = 1).values
        2  y = pandas_df["Potability"].values
```

```
330]:   1  # train test split
        2  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
        3  print("X_train",X_train.shape)
        4  print("X_test",X_test.shape)
        5  print("y_train",y_train.shape)
        6  print("y_test",y_test.shape)
```

```
X_train (2620, 9)
X_test (656, 9)
y_train (2620,)
y_test (656,)
```

Now, by training on a couple models, I see that the random forest performs better, followed by
the XGB Classifier.

```
:   1  #finding the best models with base parameters
    2  models = [
    3      ('Logistic Regression', LogisticRegression(max_iter=1000)),
    4      ('Ridge', RidgeClassifier()),
    5      ('SGD Classifier', SGDClassifier(max_iter=1000, tol=1e-3)),
    6      ('Support Vector Classifier', SVC()),
    7      ('NuSVC', NuSVC()),
    8      ('Decision Tree', DecisionTreeClassifier()),
    9      ('Gaussian NB', GaussianNB()),
   10      ('Bernoulli NB', BernoulliNB()),
   11      ('Perc', Perceptron()),
   12      ('Nearest Centroid', NearestCentroid()),
   13      ('Random Forest Classifier', RandomForestClassifier()),
   14      ('Ada Boost Classifier', AdaBoostClassifier()),
   15      ('XGB Classifier', XGBClassifier(verbosity = 0)),
   16      ('Passive Aggressive', PassiveAggressiveClassifier())
   17  ]
   18
   19  results = dict()
   20  for name, model in models:
   21      model.fit(X_train, y_train)
   22      pred = model.predict(X_test)
   23
   24
   25      score = accuracy_score(y_test, pred)
   26      mae = mean_absolute_error(y_test, pred)
   27      mse = mean_squared_error(y_test, pred)
   28      rmse = np.sqrt(mean_squared_error(y_test, pred))
   29      precision = precision_score(y_test, pred)
   30      results[name] = score
```

Here is the result figure.



## 9. HYPERPARAMETER TUNING:

This is a crucial step in optimizing the performance of machine learning models. Here, we systematically search for the best combination of hyperparameters that lead to the highest model performance on a validation dataset. By using a random search, **RandomizedSearchCV**, I found the best performer is following:

```
In [336]:   1  # Create a RandomizedSearchCV instance for RandomForestClassifier
            2  Rand_search_rf = RandomizedSearchCV(
            3      RandomForestClassifier(),              # Base classifier
            4      param_distributions=randomf_grid,      # Parameter distributions for random search
            5      cv=5,                                  # Number of cross-validation folds
            6      n_iter=20,                             # Number of parameter settings to sample
            7      verbose=True,                          # Print progress messages
            8      n_jobs=-1                              # Number of CPU cores to use (-1 for all available cores)
            9  )
           10
           11  # Fit the RandomizedSearchCV instance to the training data
           12  Rand_search_rf.fit(X_train, y_train)

           Fitting 5 folds for each of 20 candidates, totalling 100 fits

Out[336]:  RandomizedSearchCV(cv=5, estimator=RandomForestClassifier(), n_iter=20,
                   n_jobs=-1,
                   param_distributions={'max_depth': [None, 3, 5, 7, 10],
                                        'min_samples_leaf': array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 1
           3, 14, 15, 16, 17,
                   18, 19, 20, 21, 22, 23, 24]),
                                        'min_samples_split': array([ 2,  4,  6,  8, 10, 12, 14, 16, 18]),
                                        'n_estimators': array([ 10,  60, 110, 160, 210, 260, 310, 360, 410, 460, 510,
           560, 610,
                   660, 710, 760, 810, 860, 910, 960])},
                   verbose=True)
           In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
           On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [337]:   1  # Access the best parameters found by the RandomizedSearchCV
            2  best_params = Rand_search_rf.best_params_

Out[337]:  {'n_estimators': 110,
            'min_samples_split': 18,
            'min_samples_leaf': 5,
            'max_depth': 10}
```

Also, by running a grid search on the **RandomForestClassifier,** and running the accuracy test on the test set,
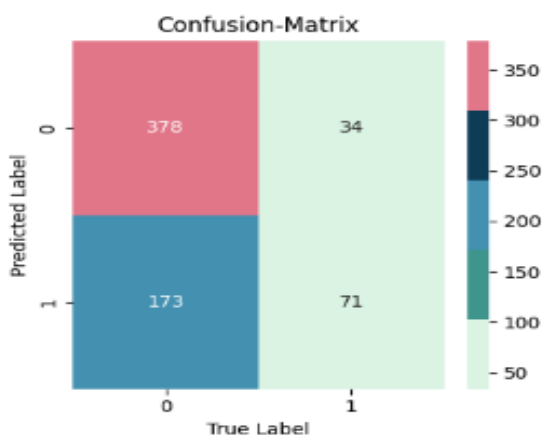
I found it working for 68.44%

```
In [345]:   1  # Calculate and print the accuracy score of the tuned model on the test set
            2  accuracy_score = gs_rf.score(X_test, y_test)
            3  print("Accuracy Score on Test Set: {} %".format(accuracy_score * 100))

Accuracy Score on Test Set: 68.4451219512195 %
```
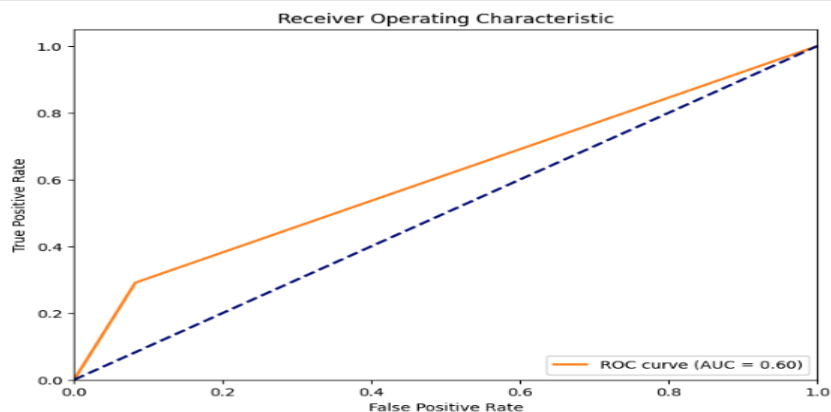
And by running the confusion matrix, the accuracy is around 68%



```
print(classification_report(y_test, y_preds))

              precision    recall  f1-score   support

         0.0       0.69      0.92      0.79       412
         1.0       0.68      0.29      0.41       244

    accuracy                           0.68       656
   macro avg       0.68      0.60      0.60       656
weighted avg       0.68      0.68      0.64       656
```

```
n [353]:    1  # Compute ROC curve and AUC
            2  fpr, tpr, thresholds = roc_curve(y_test, y_preds)
            3  roc_auc = auc(fpr, tpr)
            4
            5  # Plot ROC curve
            6  plt.figure(figsize=(8, 6))
            7  plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (AUC = %0.2f)' % roc_auc)
            8  plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
            9  plt.xlim([0.0, 1.0])
           10  plt.ylim([0.0, 1.05])
           11  plt.xlabel('False Positive Rate')
           12  plt.ylabel('True Positive Rate')
           13  plt.title('Receiver Operating Characteristic')
           14  plt.legend(loc="lower right")
           15  plt.show()
```



```
n [354]:    1  fig,ax = plt.subplots(figsize=(4,4))
            2  ax= sns.heatmap(confusion_matrix(y_test, y_preds),
```

**10. CONCLUSION:**

We have explored the application of machine learning techniques to assess water quality based on its constituent elements. Among the constructed models, the random forest exhibited the most promising outcomes. However, an inherent imbalance between the quantities of safe and unsafe samples introduced bias towards the majority class (unsafe). Consequently, various remedies for this issue were investigated, with the utilization of sample weights proving to be the most effective strategy. Nevertheless, a more in-depth exploration of this predicament is warranted to cultivate a model that offers enhanced accuracy and dependability.

In summation, the analysis of the dataset reveals the following key observations:
- ❖ The data is predominantly clean, with minimal instances of inaccurate reports that were subsequently discarded.
- ❖ The investigation did not unveil concrete evidence suggesting the redundancy of particular features or the likelihood of generating biased solutions.
- ❖ The presence of a class imbalance adversely affects the performance of classifiers.
- ❖ The application of sample weights enhances classifier accuracy, although a comprehensive inquiry remains imperative for substantiated results.

**REFERENCE:**
- ❖ Water Quality
  Kadiwal
  https://www.kaggle.com/datasets/adityakadiwal/water-potability
- ❖ Google
  www.google.com