

Water Quality Prediction using Machine Learning

Introduction

In this tutorial, we are going to implement a water quality prediction using machine learning techniques. In this technique, our model predicts that the water is safe to drink or not using some parameters like Ph value, conductivity, hardness, etc.

About dataset :-

Context

Access to safe drinking-water is essential to health, a basic human right and a component of effective policy for health protection. This is important as a health and development issue at a national, regional and local level. In some regions, it has been shown that investments in water supply and sanitation can yield a net economic benefit, since the reductions in adverse health effects and health care costs outweigh the costs of undertaking the interventions.

Content

The water_potability.csv file contains water quality metrics for 3276 different water bodies.

1. pH value:

PH 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. WHO has recommended a maximum permissible limit of pH from 6.5 to 8.5. The current investigation ranges were 6.52–6.83 which are in the range of WHO standards.

2. Hardness:

Hardness is 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.

3. Solids (*Total dissolved solids - TDS*):

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. The water with high TDS value indicates that water is highly mineralized. The Desired limit for TDS is 500 mg/l and maximum limit is 1000 mg/l which is prescribed for drinking purposes.

4. Chloramines:

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. Chlorine levels up to 4 milligrams per liter (mg/L or 4 parts per million (ppm)) are considered safe in drinking water.

5. 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. Sulfate concentration in seawater is about 2,700 milligrams per liter (mg/L). It ranges from 3 to 30 mg/L in most freshwater supplies, although much higher concentrations (1000 mg/L) are found in some geographic locations.

6. 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. Electrical conductivity (EC) actually measures the ionic process of a solution that enables it to transmit current. According to WHO standards, EC value should not exceed 400 $\mu\text{S}/\text{cm}$.

7. 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. According to the US EPA $< 2 \text{ mg/L}$ as TOC in treated / drinking water, and $< 4 \text{ mg/Lit}$ in source water which is used for treatment.

8. 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. THM levels up to 80 ppm is considered safe in drinking water.

9. 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. The mean turbidity value obtained for Wondo Genet Campus (0.98 NTU) is lower than the WHO recommended value of 5.00 NTU.

10. Potability:

Indicates if water is safe for human consumption where 1 means Potable and 0 means Not potable. (0) Water is not safe to drink and (1) Water is safe to drink.

Let's start :-

Import all the required libraries which are used to train the model or visualise the data. Then load the data set using a Pandas's function `read_csv()` and display the top five rows of the data set.

```
Data Gathering

In [7]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [8]: df = pd.read_csv("water_potability.csv")
df.head()

Out[8]:
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
0	NaN	20.898065	20.898065	1.542131	94.208054	19.21033	95.08519	2.06119	0	
1	1.78949	124.401371	188.845138	0.002136	NaN	292.861258	15.16433	30.32876	1.000004	0
2	0.269156	224.236257	188.845138	0.275684	NaN	418.88213	18.96837	96.122002	1.000004	0
3	0.187166	212.171004	220.964154	0.00133	169.88116	302.369148	18.43024	100.36174	4.020771	0
4	0.892337	131.451324	143.888338	0.548869	242.105178	345.41013	11.090278	21.867888	4.020771	0

Then perform Exploratory Data Analysis. In EDA, Firstly check the shape of the data set. Then check that there are Null values or not and you can see in the below image that ph, Sulfate, Trihalomethanes contain NULL values. Then check the information of the data set.

Exploratory Data Analysis ¶

```
In [1]: df.shape
Out[1]: (1278, 10)

In [4]: df.isnull().sum()
Out[4]:
ph                491
hardness         0
solids           0
chloramines      0
sulfate          931
conductivity      0
organic_carbon   0
trihaloethanes  182
potability       0
reliability      0
dtype: int64

In [5]: df.info()
Out[5]:
Int64Index: 1278 entries, 0 to 1277
Data columns (total 10 columns):
 #   column            non-null count  dtype
---  -
 0   ph                786 non-null    float64
 1   hardness          1278 non-null   float64
 2   solids            1278 non-null   float64
 3   chloramines       1278 non-null   float64
 4   sulfate           1249 non-null   float64
 5   conductivity      1278 non-null   float64
 6   organic_carbon    1278 non-null   float64
 7   trihaloethanes    1096 non-null   float64
 8   potability        1278 non-null   float64
 9   reliability       1278 non-null   int64
dtypes: float64(9), int64(1)
memory usage: 296.3 KB
```

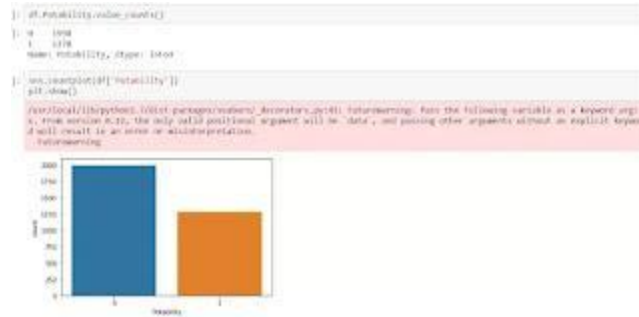
Now describe the dataset which shows the minimum value, maximum value, mean value, count, standard deviation, etc.

Then finally we handle the missing values. We filled the missing values in our features using a mean value of each feature which means we filled the mean value to handle missing data. Then again check that there are null values present or not.

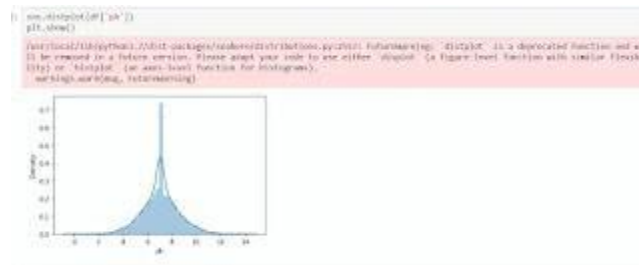
```
In [1]: df.describe()
Out[1]:
      ph  hardness  solids  chloramines  sulfate  conductivity  organic_carbon  trihaloethanes  potability  reliability
count  1278.000000  1278.000000  1278.000000  1278.000000  1249.000000  1278.000000  1278.000000  1096.000000  1278.000000  1278.000000
mean    1.024105    136.304600    1276.102200    1.122211    151.757171    438.280171    15.264678    38.292205    0.984758    0.987712
std     1.194348     21.475787     1248.176626    1.026885     81.938880     80.424003     2.384963     18.175548     0.186362    0.407508
min     0.000000     47.410000     124.940000     0.000000     15.000000     10.470000     1.000000     0.110000     0.000000     0.000000
25%     0.000000    170.000000    1000.000000     0.000000    100.000000    100.000000    10.000000     0.000000     0.000000     0.000000
50%     1.000000    140.000000    1000.000000     0.000000    150.000000    150.000000    10.000000     0.000000     0.000000     0.000000
75%     1.000000    150.000000    1250.000000     0.000000    200.000000    200.000000    10.000000     0.000000     0.000000     0.000000
max     15.000000    320.000000    4500.000000    15.000000    400.000000    700.000000    30.000000    100.000000    1.000000     1.000000

In [30]: df.fillna(df.mean(), inplace=True)
Out[30]:
ph                0
hardness          0
solids            0
chloramines       0
sulfate           0
conductivity      0
organic_carbon    0
trihaloethanes    0
potability        0
reliability       0
dtype: int64
```

Check the value counts of our target feature Potability. Then visualize the portability using a countplot function of seaborn.



Now visualize the pH value using a distplot function to check that it contains a normal distribution or not. So, you can see that it is a normal distribution.

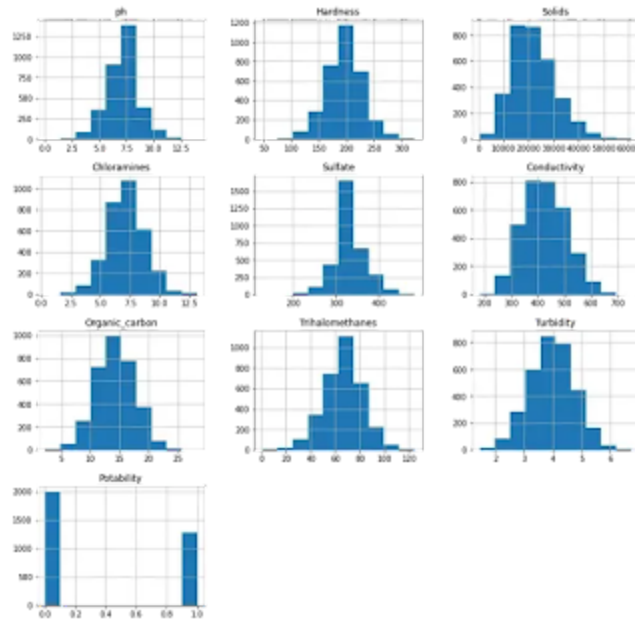


Visualize all the features of the data set as you can see below.

```

10) fig = df.hist(figsize=(14,10))
    plt.show()

```

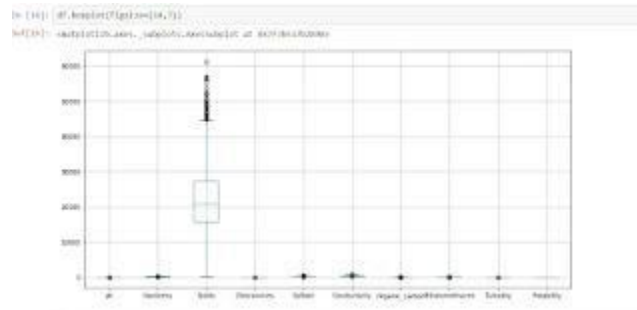


Visualize the correlation of all the features using a heat map function of seaborn. but you can see in the below heatmap that there is no correlation between any feature; it means that we can't reduce the dimension.



Now see the outlier using a boxplot function. So, you can see that the Solid feature contains outliers but we can't remove the outliers from it because if we remove the outliers from the Solid feature. So, water will be safe to drink every time. It contains an

outlier to make the water impure which means it will tell us that water is safe or not. Solid may be high to make the water unsafe to drink.



Now it's time to prepare the data set. Divide the data into the independent and dependent features. All are independent features except Potability because Potability is our dependent feature.

Split the data set into the training and testing using the `train_test_split` function which returns four data sets.

```
In [15]: X = df.drop('Potability',axis=1)
         y = df['Potability']

In [16]: from sklearn.model_selection import train_test_split
         X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.3, random_state=101,shuffle=True)
```

Now define the decision tree classifier model and train the model using the data set (`X_train, Y_train`).

Then test the model using the test data set (`X_test`).

Now it's time to evaluate the model using the accuracy score, confusion matrix and classification report. Evaluation techniques take two parameters; one is the actual data and the other one is a predicted data. And You can see that overall accuracy is 59%.


```

In [48]: from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.cross_validation import GridSearchCV, cross_val_score, cross_val_predict
from sklearn.cross_validation import train_test_split

In [49]: DecisionTreeClassifier(criterion='gini', min_samples_split=10, min_samples_leaf=10,
                               min_weight_fraction=0.05, random_state=0)

In [50]: predictions = cross_val_predict(model, X_train, y_train, cv=5)
accuracy_score(y_train, predictions)
0.92
confusion_matrix(y_train, predictions)
array([[128, 10],
       [ 9, 121]])
classification_report(y_train, predictions)
precision    recall  f1-score   support

0.00         0.00         0.00         138
1.00         0.92         0.95         131

accuracy: 0.92
macro avg: 0.50 0.50 0.50 269
weighted avg: 0.92 0.92 0.93 269

```

```

[10]: res = fit_glmfit([T5,T5T5], 150,11000,2500,80000,7.73000,107.54320,556.00000,1)
      res
fit_glmfit: 2

```

Now apply a hyper parameter tuning on the decision tree classifier using a GridSearchCV. Only three parameters to tune the model as you can see in the image below.

```
[6]: from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import grid_search_cv

# define model and parameters
model = DecisionTreeClassifier()
criteria = ['gini', 'entropy']
splitter = ['best', 'random']
min_samples_split = [3, 4, 5, 6, 10, 15, 30]

# define grid search
grid = GridSearchCV(model, criteria=criteria, min_samples=min_samples_split)
cv = RandomizedSearchCV(grid, splitter, n_repeats=10, random_state=0)
grid_search = cv.fit(X_train, y_train, scoring='accuracy', error_score='warn')
grid_search.fit(X_train, y_train)

[7]: RandomizedSearchCV(cv=DecisionTreeClassifier(), n_repeats=10, min_samples_split=3,
                        error_score='warn',
                        estimator=DecisionTreeClassifier(criterion='gini', class_weight=None,
                                                         criteria='gini', max_depth=None,
                                                         max_features=None,
                                                         max_leaf_nodes=None,
                                                         min_impurity_decrease=0.0,
                                                         min_impurity_split=None,
                                                         min_samples_leaf=1,
                                                         min_samples_split=2,
                                                         min_weight_fraction_leaf=0.0,
                                                         presort='deprecated',
                                                         random_state=None,
                                                         splitter='best'),
                        fit_params={},
                        fit_score=True,
                        fit_time=True,
                        grid_scores=True,
                        refit=True,
                        scoring='accuracy',
                        verbose=0)
```

Now see how much time the model has trained with different parameters. Also check the training and testing accuracy as you can see that the training accuracy is 90% and testing accuracy is 60% which is okay.

```
In [80]: print("Best: {grid_search_fit.best_score_:3f} using {grid_search_fit.best_params_}")
mean = grid_search_fit.cv_results_["mean_test_score"]
std = grid_search_fit.cv_results_["std_test_score"]
params = grid_search_fit.cv_results_["params"]

for mean, stdcv, param in zip(mean, std, param):
    print("%5s (%5s) with: (%5s (%5s)) with: (%5s (%5s))"

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

Best: 0.586 using {'criterion': 'gini', 'max_samples_split': 18, 'splitter': 'random'}
0.586 (0.033) with: ('criterion': 'gini', 'max_samples_split': 2, 'splitter': 'best')
0.593 (0.025) with: ('criterion': 'gini', 'max_samples_split': 2, 'splitter': 'best')
0.582 (0.030) with: ('criterion': 'gini', 'max_samples_split': 4, 'splitter': 'best')
0.585 (0.033) with: ('criterion': 'gini', 'max_samples_split': 6, 'splitter': 'random')
0.585 (0.020) with: ('criterion': 'gini', 'max_samples_split': 6, 'splitter': 'best')
0.584 (0.020) with: ('criterion': 'gini', 'max_samples_split': 6, 'splitter': 'random')
0.588 (0.021) with: ('criterion': 'gini', 'max_samples_split': 8, 'splitter': 'best')
0.577 (0.030) with: ('criterion': 'gini', 'max_samples_split': 8, 'splitter': 'random')
0.591 (0.040) with: ('criterion': 'gini', 'max_samples_split': 10, 'splitter': 'best')
0.582 (0.030) with: ('criterion': 'gini', 'max_samples_split': 10, 'splitter': 'random')
0.588 (0.035) with: ('criterion': 'gini', 'max_samples_split': 12, 'splitter': 'best')
0.588 (0.020) with: ('criterion': 'gini', 'max_samples_split': 12, 'splitter': 'random')
0.589 (0.021) with: ('criterion': 'gini', 'max_samples_split': 14, 'splitter': 'best')
0.586 (0.025) with: ('criterion': 'gini', 'max_samples_split': 14, 'splitter': 'random')
0.585 (0.032) with: ('criterion': 'entropy', 'max_samples_split': 2, 'splitter': 'best')
0.579 (0.027) with: ('criterion': 'entropy', 'max_samples_split': 2, 'splitter': 'random')
0.581 (0.030) with: ('criterion': 'entropy', 'max_samples_split': 4, 'splitter': 'best')
0.579 (0.031) with: ('criterion': 'entropy', 'max_samples_split': 6, 'splitter': 'random')
0.584 (0.030) with: ('criterion': 'entropy', 'max_samples_split': 6, 'splitter': 'best')
0.588 (0.032) with: ('criterion': 'entropy', 'max_samples_split': 6, 'splitter': 'random')
0.585 (0.025) with: ('criterion': 'entropy', 'max_samples_split': 8, 'splitter': 'best')
0.572 (0.027) with: ('criterion': 'entropy', 'max_samples_split': 8, 'splitter': 'random')
0.583 (0.030) with: ('criterion': 'entropy', 'max_samples_split': 10, 'splitter': 'best')
0.582 (0.025) with: ('criterion': 'entropy', 'max_samples_split': 10, 'splitter': 'random')
0.589 (0.033) with: ('criterion': 'entropy', 'max_samples_split': 12, 'splitter': 'best')
0.593 (0.030) with: ('criterion': 'entropy', 'max_samples_split': 12, 'splitter': 'random')
0.588 (0.029) with: ('criterion': 'entropy', 'max_samples_split': 14, 'splitter': 'best')
0.589 (0.020) with: ('criterion': 'entropy', 'max_samples_split': 14, 'splitter': 'random')
Training Score: 90.53495114549899
Testing Score: 60.808976887791
```