

Introduction to Data Science

K-Nearest Neighbors (K-NN) Project 2

Water Quality

By: Pujesh Pradhan

Table of Content

- Objective and Data Description
- Statistical Numerical and Graphical Summaries
- Algorithm Implementation and Statistical Tests
- Performance Improvement
- Conclusion

Objective and Data Description

Overview

Water is an essential part of any human being. Not all water accessible to humans are consumable. There are many factors that determine whether the water is safe for consumption. Many factors like the pH value and amount of different chemicals/minerals present in the water determines the potability of the water. Different region has different quality of water, and finding a right one can help yield an economic benefit as well as health benefits.

Objective

Predicting whether the water quality is good for consumption based on different metrics of the water.

About the Dataset

This dataset is being used from the Kaggle website and the URL to the dataset is: [Water Quality](#)

It has 10 fields and the description of each field of the dataset are:

1. **pH value:** Evaluates the acid-base balance of water. It is also the indicator of acidic or alkaline condition of water status. WHO recommends a permissible limit of pH from 6.5 to 8.5.
2. **Hardness:** Hardness defines the capacity of water to precipitate soap caused by Calcium and Magnesium. It is mainly caused by calcium and magnesium salts which are dissolved from geologic deposits.
3. **Solids (Total dissolved solids - TDS):** It is the ability to dissolve a wide range of inorganic and some organic values indicates that water is highly mineralized. Desirable limit for TDS is between 500 mg/l and 1000 mg/l which prescribed for drinking purpose.
4. **Chloramines:** Determines the amount of Chlorine and chloramine which are used as disinfectants in public water systems. Chlorine levels up to 4 milligrams per liter (mg/l or 4 parts per million (ppm)) are considered safe for drinking water.
5. **Sulfate:** Indicates the amount of Sulfates in the water. They are naturally occurring substances that are found in minerals, soil, and rocks.
6. **Conductivity:** Indicates the ionic conductivity of the water. Electrical conductivity (EC) actually measures the ionic process of a solution that enables it to transmit current. Pure water is not a good conductor of electric current rather it is a good insulator. The amount of dissolved solids in water determines the electrical conductivity. According to WHO standards, EC value should not exceed 400 µS/cm.
7. **Organic_carbon:** Total Organic Carbon (TOC) measures the total amount of carbon in organic compounds in pure water. It comes from decaying natural organic matter (NOM) as well as synthetic sources.
8. **Trihalomethanes:** 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:** 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. It depends on the quantity of solid matter present in the suspended state.
10. **Potability:** Indicates if water is safe for human consumption where 1 means Potable and 0 means Not potable.

```
In [1]: #All the necessary packages are imported here.
import numpy as np
import pandas as pd
from pandas import scatter_matrix
import seaborn as sns
from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score
from sklearn.model_selection import cross_val_score, train_test_split
from sklearn.preprocessing import MinMaxScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestRegressor
from math import sqrt
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: #Describing the size, shape and dimension of the dataset.
water_metrics = pd.read_csv("water_potability.csv")
print("The total size of the dataset is (%d)." % water_metrics.size)
print("There are (%d) number of records and (%d) number of columns." % (water_metrics.shape[0], water_metrics.shape[1]))
```

```
The total size of the dataset is 32760 bytes.
The dataset has 3276 number of rows and 10 number of fields with a (3276, 10) shape.
The dataset has a 2 dimension structure.
```

```
In [3]: print("The top 5 records of the dataset are: \n").format(water_metrics.head()))
The top 5 records of the dataset are:
   ph  hardness  Solids  Chloramines  Sulfate  Conductivity
0    6.937978  204.18009  207.91100  7.000212  368.516441  564.309654
1    3.716080  129.422921  186.30105  0.057858  6.635246      NaN  592.485359
2    8.099124  220.6.2204.09.54.26.59  199.0.54.17.41  9.275884  59.418.06213
3    8.516764  214.3.73.39.4  202.0.41.74.41  8.0.59.32  356.88.61.36  363.266516
4    9.092233  181.1.01.059  179.78.98.63.99  5.64.66.00  310.13.57.98  398.410813
```

```
Here certain fields have missing values that are filled by NaN values.
```

```
In [4]: print("The datatype of the dataset are: \n").format(water_metrics.dtypes))
The datatype of the dataset are:
ph          float64
Hardness     float64
Solids      float64
Chloramines  float64
Sulfate      float64
Conductivity float64
Organic_carbon float64
Trihalomethanes float64
Turbidity    float64
Potability   int64
dtype: object
```

```
In [5]: print(water_metrics.describe().T)
```

```
      count      mean        std       min      25%      50%      max
ph      3276.0  7.036752  8.052969  14.000000  5.639330  6.089723
Hardness  3276.0  197.191839  216.410170  317.338124
Solids   3276.0  2093.513645  2198.500000  5646.242413
Chloramines  3276.0  7.134338  8.114887  13.127000
Sulfate   3276.0  322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  3276.0  426.526409  80.712572  201.619737  366.680307
Organic_carbon  3276.0  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  3114.0  66.396293  16.175000  0.780346  2.200000  12.065801
Turbidity   3276.0  3.955028  4.500320  6.737900  1.450000  3.439711
Potability   3276.0  0.390110  0.487849  0.000000  0.000000  0.000000
      50%      75%      max
ph      7.027297  8.052969  14.000000
Hardness  197.191839  216.410170  317.338124
Solids   2093.513645  2198.500000  5646.242413
Chloramines  7.134338  8.114887  13.127000
Sulfate   322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  426.526409  80.712572  201.619737  366.680307
Organic_carbon  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  66.396293  16.077109  8.577013  55.532664
Turbidity   3.955028  4.500320  6.737900
Potability   0.390110  0.487849  0.000000
```

```
From the above statistics, we can see that there are certain values missing for ph, sulphate and trihalomethanes. The metrics min ranges widely from each other, thus scalability/normalization of the data is required. Few of the metrics looks to display skewness in the data.
```

```
We will be removing all instances that has NaN values in it.
```

```
In [6]: #Deleting all the rows with NaN values
water_metrics = water_metrics.dropna()
```

```
In [7]: print("Now, the number of rows has been reduced to () rows and the new shape is ()".format(water_metrics.shape))
print("The new data looks like this: \n").format(water_metrics.describe().T)
```

```
Now, the number of rows has been reduced to 2011 rows and the new shape is (2011, 10)
The new data looks like this:
   ph  hardness  Solids  Chloramines  Sulfate  Conductivity
0    6.937978  204.18009  207.91100  7.000212  368.516441  564.309654
1    3.716080  129.422921  186.30105  0.057858  6.635246      NaN  592.485359
2    8.099124  220.6.2204.09.54.26.59  199.0.54.17.41  9.275884  59.418.06213
3    8.516764  214.3.73.39.4  202.0.41.74.41  8.0.59.32  356.88.61.36  363.266516
4    9.092233  181.1.01.059  179.78.98.63.99  5.64.66.00  310.13.57.98  398.410813
      count      mean        std       min      25%      50%      max
ph      2011.0  7.085993  8.052969  14.000000  5.639330  6.089723
Hardness  2011.0  195.204.180.09.54.26.59  207.91100  7.000212  368.516441  564.309654
Solids   2011.0  2093.513645  2198.500000  5646.242413
Chloramines  2011.0  7.134338  8.114887  13.127000
Sulfate   2011.0  322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  2011.0  426.526409  80.712572  201.619737  366.680307
Organic_carbon  2011.0  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  2011.0  66.396293  16.077109  8.577013  55.532664
Turbidity   2011.0  3.955028  4.500320  6.737900
Potability   2011.0  0.390110  0.487849  0.000000  0.000000  0.000000
      50%      75%      max
ph      7.027297  8.052969  14.000000
Hardness  197.191839  216.410170  317.338124
Solids   2093.513645  2198.500000  5646.242413
Chloramines  7.134338  8.114887  13.127000
Sulfate   322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  426.526409  80.712572  201.619737  366.680307
Organic_carbon  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  66.396293  16.077109  8.577013  55.532664
Turbidity   3.955028  4.500320  6.737900
Potability   0.390110  0.487849  0.000000
```

```
From the above statistics, we can see that there are certain values missing for ph, sulphate and trihalomethanes. The metrics min ranges widely from each other, thus scalability/normalization of the data is required. Few of the metrics looks to display skewness in the data.
```

```
We will be removing all instances that has NaN values in it.
```

```
In [6]: #Deleting all the rows with NaN values
water_metrics = water_metrics.dropna()
```

```
In [7]: print("Now, the number of rows has been reduced to () rows and the new shape is ()".format(water_metrics.shape))
print("The new data looks like this: \n").format(water_metrics.describe().T)
```

```
Now, the number of rows has been reduced to 2011 rows and the new shape is (2011, 10)
The new data looks like this:
   ph  hardness  Solids  Chloramines  Sulfate  Conductivity
0    6.937978  204.18009  207.91100  7.000212  368.516441  564.309654
1    3.716080  129.422921  186.30105  0.057858  6.635246      NaN  592.485359
2    8.099124  220.6.2204.09.54.26.59  199.0.54.17.41  9.275884  59.418.06213
3    8.516764  214.3.73.39.4  202.0.41.74.41  8.0.59.32  356.88.61.36  363.266516
4    9.092233  181.1.01.059  179.78.98.63.99  5.64.66.00  310.13.57.98  398.410813
      count      mean        std       min      25%      50%      max
ph      2011.0  7.085993  8.052969  14.000000  5.639330  6.089723
Hardness  2011.0  195.204.180.09.54.26.59  207.91100  7.000212  368.516441  564.309654
Solids   2011.0  2093.513645  2198.500000  5646.242413
Chloramines  2011.0  7.134338  8.114887  13.127000
Sulfate   2011.0  322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  2011.0  426.526409  80.712572  201.619737  366.680307
Organic_carbon  2011.0  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  2011.0  66.396293  16.077109  8.577013  55.532664
Turbidity   2011.0  3.955028  4.500320  6.737900
Potability   2011.0  0.390110  0.487849  0.000000  0.000000  0.000000
      50%      75%      max
ph      7.027297  8.052969  14.000000
Hardness  197.191839  216.410170  317.338124
Solids   2093.513645  2198.500000  5646.242413
Chloramines  7.134338  8.114887  13.127000
Sulfate   322.232177  41.584820  41.205172  129.000000  307.632511
Conductivity  426.526409  80.712572  201.619737  366.680307
Organic_carbon  14.218398  16.576562  80.651300  102.337800  364.221414
Trihalomethanes  66.396293  16.077109  8.577013  55.532664
Turbidity   3.955028  4.500320  6.737900
Potability   0.390110  0.487849  0.000000
```

```
From the above statistics, we can see that there are certain values missing for ph, sulphate and trihalomethanes. The metrics min ranges widely from each other, thus scalability/normalization of the data is required. Few of the metrics looks to display skewness in the data.
```

```
We will be removing all instances that has NaN values in it.
```

```
In [6]: #Deleting all the rows with NaN values
water_metrics = water_metrics.dropna()
```

```
In [7]: print("Now, the number of rows has been reduced to () rows and the new shape is ()".format(water_metrics.shape))
print("The new data looks like this: \n").format(water_metrics.describe().T)
```

```
Now, the number of rows has been reduced to 2011 rows and the new shape is (2011, 10)
The new data looks like this:
   ph  hardness  Solids  Chloramines  Turbidity  Potability
0    6.937978  204.18009  207.91100  7.000212  368.516441  564.309654
1    3.716080  129.422921  186.30105  0.057858  6.635246      NaN  592.485359
2    8.099124  220.6.2204.09.54.26.59  199.0.54.17.41  9.275884  59.418.06213
3    8.516764  214.3.73.39.4  202.0.41.74.41  8.0.59.32  356.88.61.36  363.266516
4    9.092233 
```

