

# A Study of Water Potability

An attempt on understanding what makes water potable based on 10 different numerical features. The water\_potability.csv file contains water quality metrics for 3276 different water bodies and its source can be found here: <https://www.kaggle.com/datasets/ashishkashish/water-potability?datasetid=1232407&entry=code&view=code>

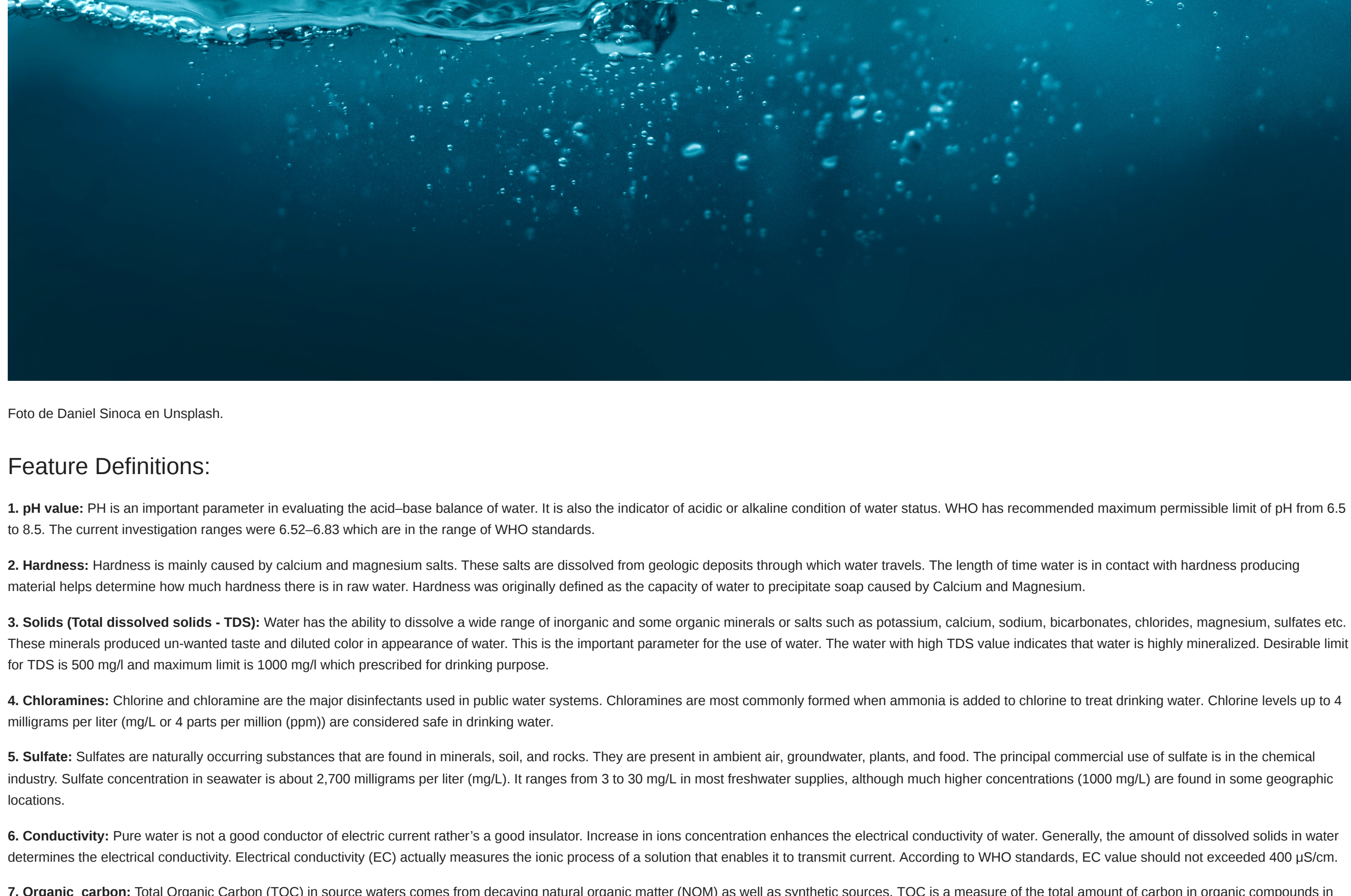


Foto: De Daniel Sica on Unplash.

## Feature Definitions:

- 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 maximum permissible limit of pH from 6.5 to 9.5. The current investigation ranges were 6.5 to 6.83 which are in the range of WHO standards.
- 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.
- 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 produce un-wanted taste and diluted color in 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. Desirable limit for TDS is 500 mg/l and maximum limit is 1000 mg/l which prescribed for drinking purpose.
- 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 a part per million (ppm)) are considered safe for 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. 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 geologic locations.
- Conductivity:** Pure water is not a good conductor of electric current rather 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  $\mu$ S/cm.
- Organic carbon:** Total Organic Carbon (TOC) in source water 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 US EPA <2 mg/L as TOC in treated / drinking water, and <4 mg/L in its source water which is use for treatment.
- 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.
- 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 water 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.
- Potability:** Indicates if water is safe for human consumption where 1 means Potable and 0 means not potable.

## 0.1 Import Libraries & Packages

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
plt.rcParams["figure.figsize"] = (20,10)
import warnings as wps

from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.pipeline import Pipeline

In [2]: # set seed
seed = 42
np.random.seed(seed)

# Import and show data
file = 'data/water_potability.csv'
data = pd.read_csv(file)
data.head()
```

```
Out[2]:
```

	pH	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
0	NAN	204.960455	10780.318981	7.302212	368.516441	564.309854	10.379783	65.950570	2.963115	0
1	3.710080	129.424281	20350.057858	6.036346	NAN	592.880339	15.180013	56.529676	4.500565	0
2	8.099524	224.236259	19950.541732	9.275884	NAN	418.686337	18.868637	66.420093	3.055034	0
3	8.316766	214.373394	22018.417441	8.059332	366.886136	363.265513	18.436524	100.341674	4.628771	0
4	8.092223	181.101509	17978.886329	6.546600	310.129738	356.412613	11.558759	31.907963	4.075075	0

## 1.0 Exploratory Data Analysis

### 1.1 Quick Descriptive Statistics

We will be exploring general characteristics for the data. Things like type and name of each variable, missing values, dimensions of the data, and more.

```
In [3]: # Show data types, nulls, features, and dimensions.
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3276 entries, 0 to 3275
Data columns (total 10 columns):
 #   Column                Non-Null Count  Dtype  ---
 0   pH                    2785 non-null   float64
 1   Hardness              3276 non-null   float64
 2   Solids                3276 non-null   float64
 3   Chloramines           3276 non-null   float64
 4   Sulfate               2495 non-null   float64
 5   Conductivity          3276 non-null   float64
 6   Organic_carbon        3276 non-null   float64
 7   Trihalomethanes       2114 non-null   float64
 8   Turbidity             3276 non-null   float64
 9   Potability            3276 non-null   int64
dtypes: float64(9), int64(1)
memory usage: 258.1 KB

In [4]: data.describe()
```

```
Out[4]:
```

	pH	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
count	2785.000000	3276.000000	3276.000000	2495.000000	3276.000000	3276.000000	3114.000000	3276.000000	3276.000000	3276.000000
mean	7.907096	196.309696	22014.002526	7.122277	333.775777	426.205111	14.246870	66.398038	3.960786	0.390120
std	1.504330	32.877951	8769.578509	1.503095	41.415040	80.824004	3.306162	16.175009	0.780382	0.487849
min	0.000000	47.432000	320.942611	0.362000	129.000000	181.483754	2.200000	0.738000	0.000000	0.000000
25%	6.000000	176.650638	15566.690297	6.127421	307.694981	365.734414	12.059001	55.844036	3.439711	0.000000
50%	7.090762	196.967627	22027.833607	7.130209	333.073546	421.884968	14.218308	66.624468	3.955028	0.000000
75%	8.062006	214.667456	27332.762127	8.114887	359.850170	481.782034	16.751502	77.321473	4.500200	1.000000
max	14.000000	325.124500	61227.199008	13.127000	481.209842	763.342620	28.300000	124.000000	6.799000	1.000000

```
In [5]: # Percentage of missing values.
print('Percentage of missing values for each variable:')
data.isnull().sum()[data.columns[0]]

Percentage of missing values for each variable:
pH                0.246179
Hardness          0.000000
Solids            0.000000
Chloramines       0.000000
Sulfate           0.228482
Conductivity      0.000000
Organic_carbon    0.000000
Trihalomethanes  0.348451
Turbidity         0.000000
Potability        0.000000
dtype: float64

In [6]: We can see that each variable has a different range of values when describing the potability of a body of water. This is important because it will determine the necessary transformations on the data to create an accurate predictive model. It is important to note that the variables: pH, sulfate, and trihalomethanes present a significant amount of missing values that will need to be addressed.
```

### 1.2 Data Distributions

By visualizing the distributions of each independent variables we will be able to have a more accurate understanding of the relationship between the respective dependent variable. We will make two plots: one for the entire data, and the other for each variable conditioned on its label.

```
In [6]: # Plot Features Distributions
fig, axs = plt.subplots(3,3)
for i in range(data.shape[1]-1):
    plt.subplot(3,3,i+1)
    sns.kdeplot(data[data.columns[i]])

    Distribution of Each Feature
```

We can see that most variables are already normally distributed with solids having a slight skew to the right. This is important to know when deciding which predictive model to use as some of them make normality of independent variables an assumption to work. We will now visualize the relationship between the variables.

```
In [7]: # Plot Features Distributions
fig, axs = plt.subplots(3,3)
for i in range(data.shape[1]-1):
    plt.subplot(3,3,i+1)
    pot_var, non_pot_var = df_pot[data.columns[i]], df_non_pot[data.columns[i]]
    sns.kdeplot(data[pot_var, 'Potability'])
    sns.kdeplot(data[non_pot_var, 'Potability'])

    Distribution by Potability (Red: Non-potable, Green: Potable)
```

We do this to gain some insights between the dependent and independent variables. We can see how some variables conditional distribution remains the same regardless of class, but others change slightly like pH, sulfate, or hardness. Overall, the differences seem to be small which might represent a challenge when trying to predict new data.

Another way to gain insight on the variables relationship and a good method to check for multicollinearity is to do a correlation plot.

```
In [8]: #Feature correlation
corr_heat = sns.heatmap(data.corr(),annot=True)
```

In it we can see the correlation between all the variables. In this case we can see that there are no strong relationships among any of the features. The good side of this is that we don't have to worry about the model being affected by linear relationships between independent variables, the bad is that it might be hard predicting the test data class due to the lack of relationship with the dependent variable.

### 1.3 Outlier Detection

We will use boxplots for each variable to have a visual intuition for each variables deviations. This will help gain insight on the necessary transformations for modeling.

```
In [9]: # Outlier detection
fig, axs = plt.subplots(3,3)
for i in range(9):
    plt.subplot(3,3,i+1)
    sns.boxplot(x=data[data.columns[i+1]])

    Distribution of Each Feature
```

```
In [10]: # Outlier detection
fig, axs = plt.subplots(3,3)
for i in range(9):
    plt.subplot(3,3,i+1)
    sns.boxplot(x=data[data.columns[i+1]])

    Distribution of Each Feature
```

We can see that a lot of variables present a significant number of outliers. We will either need to make transformations for each variable or select a model robust to outliers.

### 1.4 Check Data for Imbalance & Split Into Train/Test

```
In [11]: # Distribution of dependent variable.
data['Potability'].value_counts()

Out[11]:
0    1998
1     128
Name: Potability, dtype: int64

We can see how there exists a slight imbalance between the 2 classes in our data, this might be important to take into account when building our model.

We will now split our data to have a training and test set. We will do this by first separating the data points with missing values and splitting the rest. We will then combine the newly formed training set with the data points with missing values. This will prevent data leakage into the test set.

In [12]: # Separate missing data from complete data for splitting before preprocessing.
complete_data = data.dropna()
non_pot = data[data['Potability'] == 0]
pot = data[data['Potability'] == 1]

In [13]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [14]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [15]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [16]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [17]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [18]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [19]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [20]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [21]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [22]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [23]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [24]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [25]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [26]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [27]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [28]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [29]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [30]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [31]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [32]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [33]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [34]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [35]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [36]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [37]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [38]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [39]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [40]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [41]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [42]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [43]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [44]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [45]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [46]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [47]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [48]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [49]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [50]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [51]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [52]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [53]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [54]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [55]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [56]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [57]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [58]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [59]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [60]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [61]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [62]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [63]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [64]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [65]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [66]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [67]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [68]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [69]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [70]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [71]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [72]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [73]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [74]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [75]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [76]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [77]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [78]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [79]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [80]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [81]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [82]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [83]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [84]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [85]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [86]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [87]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [88]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [89]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [90]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [91]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [92]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [93]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [94]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [95]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [96]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [97]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [98]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [99]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [100]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [101]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [102]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [103]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [104]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [105]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [106]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [107]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [108]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [109]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [110]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)

In [111]: # Split data
X, y = complete_data.loc[:, complete_data.columns != 'Potability'], complete_data['Potability']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=seed)
```