

```
In [37]: import pandas as pd
import numpy as np
```

```
In [38]: data=pd.read_csv("C:/Users/DELL/Desktop/water_potability.csv")
```

```
In [39]: data.head()
```

| | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|----------|------------|--------------|-------------|------------|--------------|----------------|-----------------|-----------|------------|
| 0 | NaN | 204.890455 | 20791.318981 | 7.300212 | 368.516441 | 564.308654 | 10.379783 | 86.990970 | 2.963135 | 0 |
| 1 | 3.716080 | 129.422921 | 18630.057858 | 6.635246 | NaN | 592.885359 | 15.180013 | 56.329076 | 4.500656 | 0 |
| 2 | 8.099124 | 224.236259 | 19909.541732 | 9.275884 | NaN | 418.606213 | 16.868637 | 66.420093 | 3.055934 | 0 |
| 3 | 8.316766 | 214.373394 | 22018.417441 | 8.059332 | 356.886136 | 363.266516 | 18.436524 | 100.341674 | 4.628771 | 0 |
| 4 | 9.092223 | 181.101509 | 17978.986339 | 6.546600 | 310.135738 | 398.410813 | 11.558279 | 31.997993 | 4.075075 | 0 |

```
In [40]: data.info
```

| | | | | | | | |
|----------|---------------------------------|---------------------------|----------------|-----------------|-------------|------------|---|
| Out[40]: | <bound method DataFrame.info of | ph | Hardness | Solids | Chloramines | Sulfate | \ |
| 0 | 0 | NaN | 204.890455 | 20791.318981 | 7.300212 | 368.516441 | |
| 1 | 1 | 3.716080 | 129.422921 | 18630.057858 | 6.635246 | NaN | |
| 2 | 2 | 8.099124 | 224.236259 | 19909.541732 | 9.275884 | NaN | |
| 3 | 3 | 8.316766 | 214.373394 | 22018.417441 | 8.059332 | 356.886136 | |
| 4 | 4 | 9.092223 | 181.101509 | 17978.986339 | 6.546600 | 310.135738 | |
| ... | ... | ... | ... | ... | ... | ... | |
| 3271 | 3271 | 4.668102 | 193.681735 | 47580.991603 | 7.166639 | 359.948574 | |
| 3272 | 3272 | 7.808856 | 193.553212 | 17329.802160 | 8.061362 | NaN | |
| 3273 | 3273 | 9.419510 | 175.762646 | 33155.578218 | 7.350233 | NaN | |
| 3274 | 3274 | 5.126763 | 230.603758 | 11983.869376 | 6.303357 | NaN | |
| 3275 | 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509306 | NaN | |
| | | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability | |
| 0 | 0 | 564.308654 | 10.379783 | 86.990970 | 2.963135 | 0 | |
| 1 | 1 | 592.885359 | 15.180013 | 56.329076 | 4.500656 | 0 | |
| 2 | 2 | 418.606213 | 16.868637 | 66.420093 | 3.055934 | 0 | |
| 3 | 3 | 363.266516 | 18.436524 | 100.341674 | 4.628771 | 0 | |
| 4 | 4 | 398.410813 | 11.558279 | 31.997993 | 4.075075 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | |
| 3271 | 3271 | 526.424171 | 13.894419 | 66.687695 | 4.435821 | 1 | |
| 3272 | 3272 | 392.449580 | 19.903225 | NaN | 2.798243 | 1 | |
| 3273 | 3273 | 432.044783 | 11.039070 | 69.845400 | 3.298875 | 1 | |
| 3274 | 3274 | 402.883113 | 11.168946 | 77.488213 | 4.708658 | 1 | |
| 3275 | 3275 | 327.459760 | 16.140368 | 78.698446 | 2.309149 | 1 | |
| | | [3276 rows x 10 columns]> | | | | | |

```
In [41]: data.dtypes
```

| | | |
|----------|-----------------|---------|
| Out[41]: | ph | float64 |
| | Hardness | float64 |
| | Solids | float64 |
| | Chloramines | float64 |
| | Sulfate | float64 |
| | Conductivity | float64 |
| | Organic_carbon | float64 |
| | Trihalomethanes | float64 |
| | Turbidity | float64 |
| | Potability | int64 |
| | dtype: object | |

```
In [42]: missing_values = data.isnull().sum()
```

```
In [43]: print("missing values:\n",missing_values)

missing values:
  ph      491
Hardness    0
Solids      0
Chloramines  0
Sulfate    781
Conductivity  0
Organic_carbon  0
Trihalomethanes 162
Turbidity    0
Potability    0
dtype: int64
```

```
In [44]: #2 Define outlier and its importance in data analysis
# Outliers refer to observations that significantly deviate from other data points in a dataset. They can affect statistical analyses, leading to in
```

```
In [45]: #2b Importance in data analysis
# Outliers can impact statistical significance, model performance, and data quality. They can affect mean and standard deviation, influence machine
```

```
In [46]: #Detection of of outliers in datasets
# Box and scatter plots visually represent data distribution and identify potential outliers. Z-score measures standard deviations from the mean, wi
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```