# Water quality analysis

## PHASE:3 PROJECT

TEAM MEMBERS

D.SHREE THOVARTHENE

G.NILOFAR

E.PREETHI RAJALAKSHMI

K.MAHALAKSHMI

- PREPROCESSING THE DATASET

- PERFORMING ANALYSIS AND VISUALIZATION USING IBM COGNOS
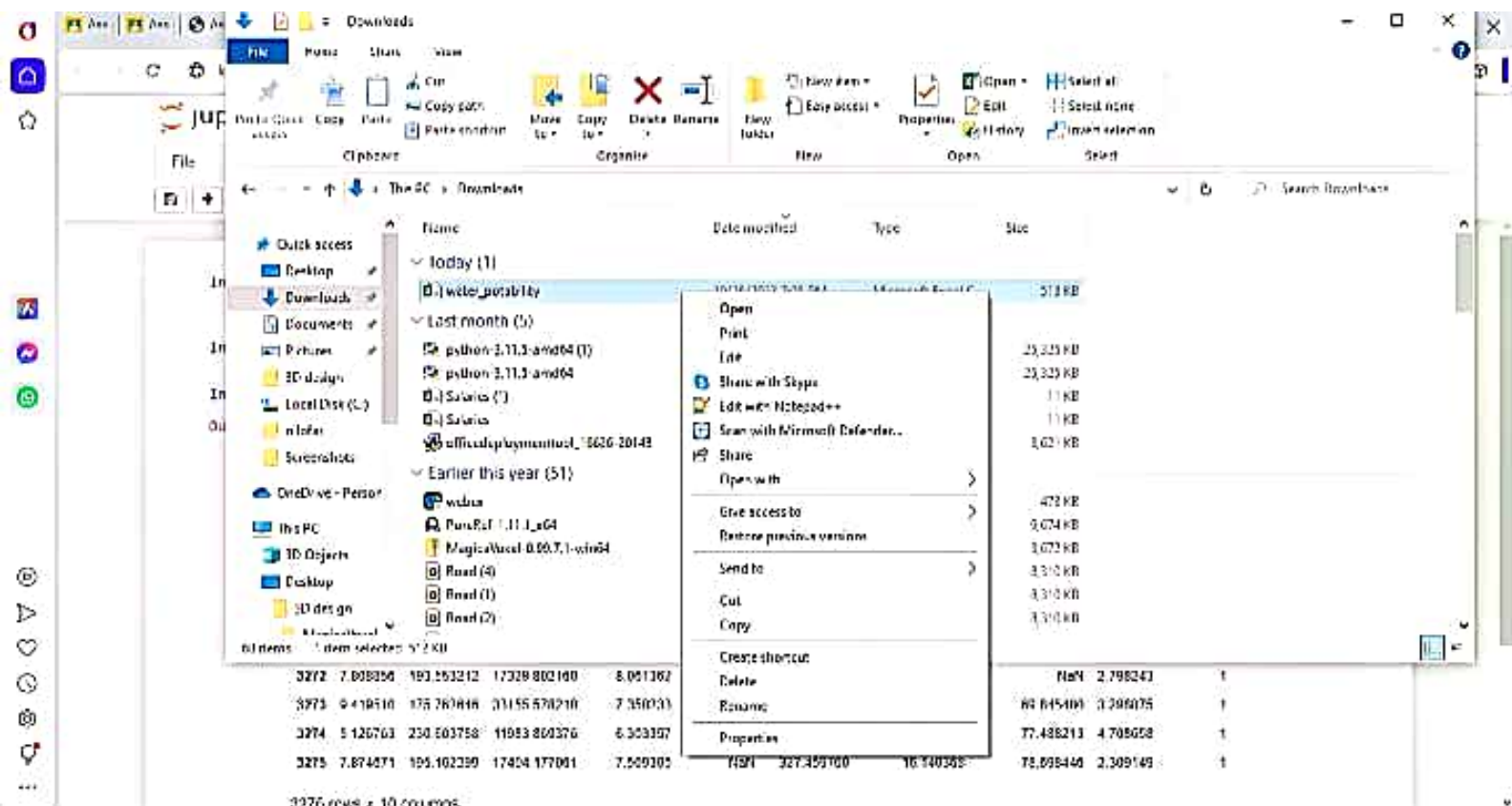
# Preprocessing the dataset:

## Definition:

Preprocessing a dataset is a crucial step in data analysis and machine learning. It involves cleaning and transforming the data to make it suitable for analysis or model training.
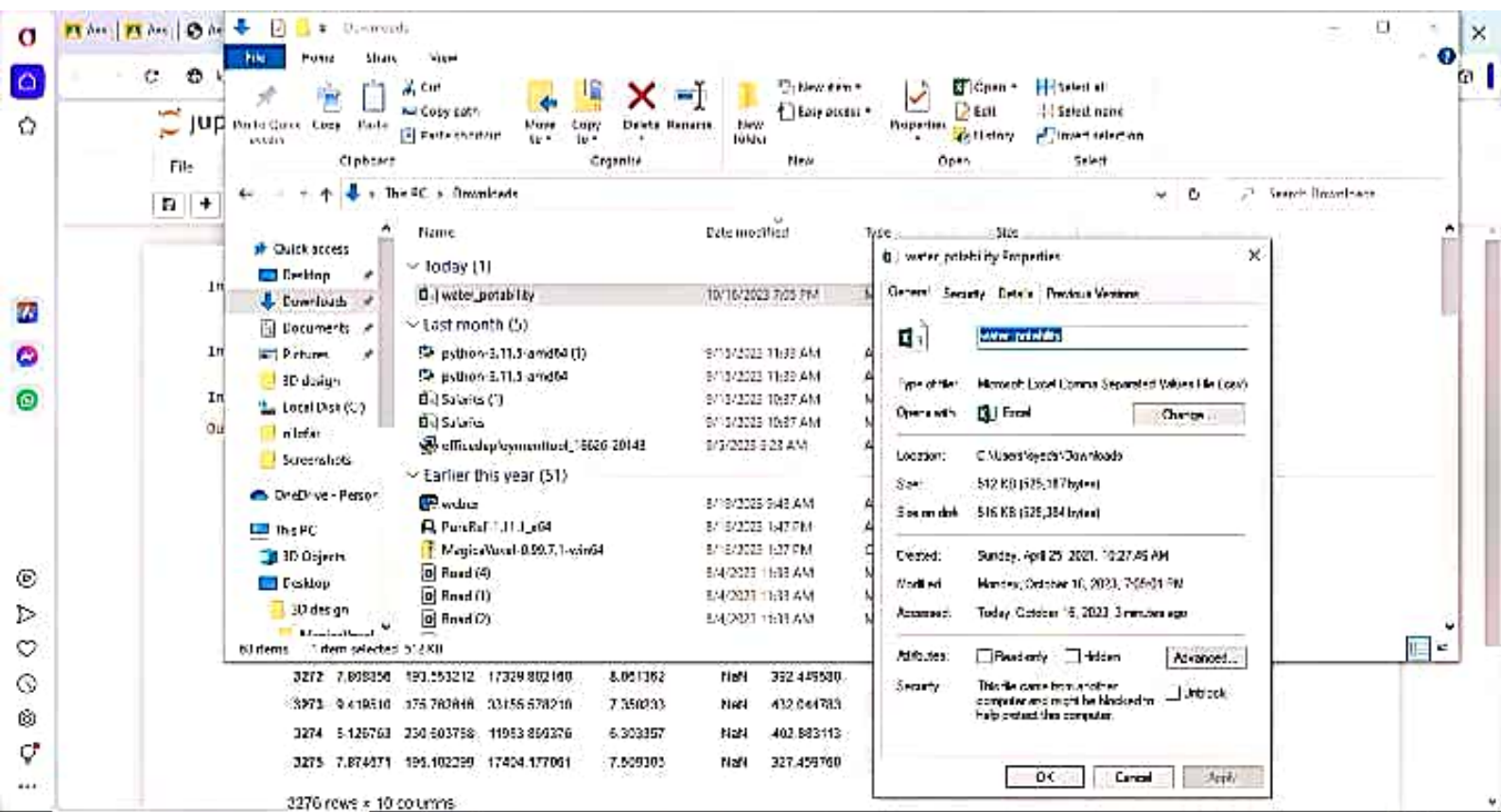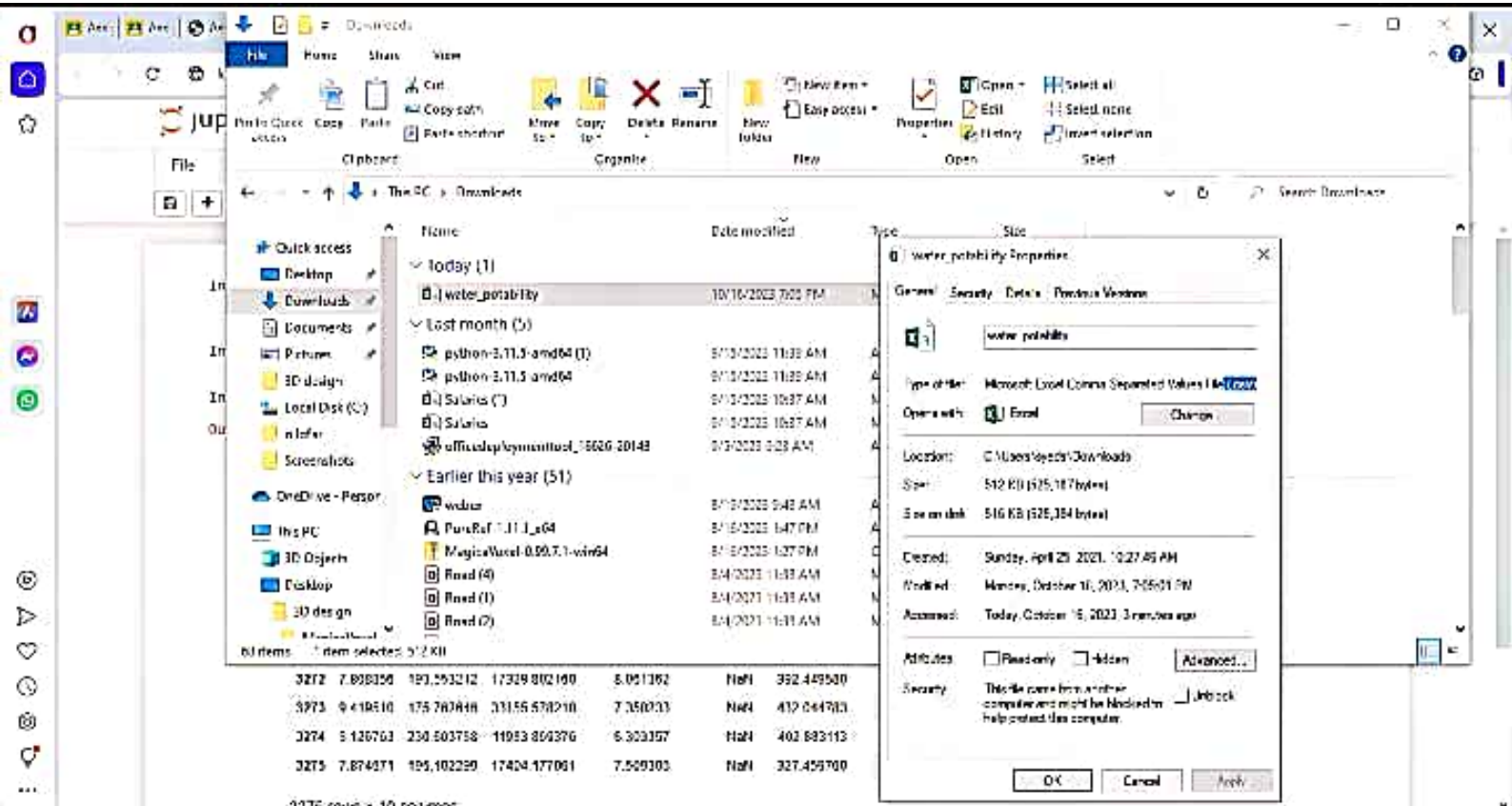
## Data cleaning:

- Data cleaning is an essential part of preprocessing your dataset to ensure its quality and reliability.

## Handling Missing Data:

- Identify and analyze missing values in your dataset.
- Decide whether to impute missing values, remove rows or columns with missing data, or use other techniques like interpolation.

| 3272 | 7.808856 | 193.553212 | 17328.802160 | 8.061362 |  | NaN | 2.798243 | t |
| 3273 | 9.419510 | 175.760816 | 33155.578210 | 7.350033 |  | 89.845400 | 3.298075 | t |
| 3274 | 5.126763 | 230.603758 | 11983.869376 | 6.303367 |  | 77.488213 | 4.708658 | t |
| 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509302 | NaN | 327.459760 | 16.140368 | 78.698446 | 2.309149 | t |

3276 rows × 10 columns

**File Explorer - Downloads**

Navigation pane:
- Quick access
- Desktop
- Downloads
- Documents
- Pictures
- 3D design
- Local Disk (C:)
- n infer
- Screenshots
- OneDrive - Person
- This PC
- 3D Objects
- Desktop
- 3D design

| Name | Date modified | Type | Size |
|---|---|---|---|
| **Today (1)** | | | |
| water_potability | 10/16/2023 7:05 PM | | |
| **Last month (5)** | | | |
| python-3.11.5-amd64 (1) | 9/13/2023 11:33 AM | | |
| python-3.11.5-amd64 | 9/13/2023 11:33 AM | | |
| Salaries (1) | 9/13/2023 10:37 AM | | |
| Salaries | 9/13/2023 10:37 AM | | |
| officialsymantico1_16626 20148 | 9/3/2023 9:23 AM | | |
| **Earlier this year (51)** | | | |
| wakaa | 9/19/2023 9:43 AM | | |
| PureRef-1.11.1_x64 | 9/13/2023 1:47 PM | | |
| MagicaVoxel-0.99.7.1-win64 | 9/13/2023 1:27 PM | | |
| Read (4) | 9/4/2023 11:33 AM | | |
| Read (1) | 9/4/2023 11:33 AM | | |
| Read (2) | 9/4/2023 11:33 AM | | |

60 items    1 item selected 512 KB

**water_potability Properties**

General | Security | Details | Previous Versions

water_potability

| | |
|---|---|
| Type of file: | Microsoft Excel Comma Separated Values File (.csv) |
| Opens with: | Excel    Change... |
| Location: | C:\Users\iejeda\Downloads |
| Size: | 512 KB (525,187 bytes) |
| Size on disk: | 516 KB (525,384 bytes) |
| Created: | Sunday, April 25, 2021, 10:27:46 AM |
| Modified: | Monday, October 16, 2023, 7:05:01 PM |
| Accessed: | Today, October 16, 2023, 3 minutes ago |

Attributes:  ☐ Read-only   ☐ Hidden   Advanced...

Security: This file came from another computer and might be blocked to help protect this computer.   ☐ Unblock

OK | Cancel | Apply

| | | | | | | |
|---|---|---|---|---|---|---|
| 3272 | 7.808356 | 190.553212 | 17329.802168 | 8.061362 | NaN | 392.449580 |
| 3273 | 9.419510 | 175.762848 | 33155.578216 | 7.350233 | NaN | 432.044783 |
| 3274 | 5.126763 | 230.603758 | 11983.869376 | 6.303357 | NaN | 402.883113 |
| 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509303 | NaN | 327.459760 |

3276 rows × 10 columns

```
In [20]: import numpy as np
         import pandas as pd
```

```
In [23]: d = pd.read_csv("water_potability.csv")
```

```
In [24]: d
```

Out[24]:

| | 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 | 18630.057860 | 6.635245 | 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 |
| ... | | | | | | | | | | |
| 3271 | 4.668102 | 193.681735 | 47580.991603 | 7.166639 | 359.948574 | 526.424171 | 13.894419 | 66.687695 | 4.435821 | 1 |
| 3272 | 7.808856 | 193.553212 | 17329.802160 | 8.061362 | NaN | 392.449580 | 19.903225 | NaN | 2.798243 | 1 |
| 3273 | 9.419510 | 175.762646 | 33155.578218 | 7.350233 | NaN | 432.044783 | 11.039070 | 69.845400 | 3.298875 | 1 |
| 3274 | 5.126763 | 230.603758 | 11983.869376 | 6.303357 | NaN | 402.883113 | 11.168946 | 77.488213 | 4.708658 | 1 |
| 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509306 | NaN | 327.459760 | 16.140368 | 78.698446 | 2.309149 | 1 |

3276 rows × 10 columns

```
3274  5.126763  230.803750  11953.889376  6.303357  NaN  402.883113  11.168045  77.488211  4.708658   1
3275  7.874671  195.102298  17454.177091  7.568305  NaN  317.450760  16.140368  78.680046  2.309149   1
```

3276 rows × 10 columns

In [25]: d.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  3114 non-null   float64
 8   Turbidity        3276 non-null   float64
 9   Potability       3276 non-null   int64
dtypes: float64(9), int64(1)
memory usage: 256.1 KB
```

In [ ]:

```
7  Trihalomethanes  3114 non-null   float64
8  Turbidity        3276 non-null   float64
9  Potability       3276 non-null   int64
dtypes: float64(9), int64(1)
memory usage: 256.1 KB
```

In [26]: d.dropna()

Out[26]:

|  | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 0.314766 | 214.373394 | 22018.417441 | 8.059332 | 356.886136 | 363.266516 | 18.436524 | 100.341674 | 4.628771 | 0 |
| 4 | 9.092223 | 181.101509 | 17978.986332 | 6.546600 | 310.135738 | 398.410813 | 11.558279 | 31.997993 | 4.075075 | 0 |
| 5 | 5.584087 | 188.313324 | 28748.687739 | 7.544869 | 326.678363 | 280.467916 | 8.399735 | 54.917862 | 2.559708 | 0 |
| 6 | 10.223862 | 248.071735 | 28749.716544 | 7.513408 | 393.663396 | 283.651634 | 13.789695 | 84.603556 | 2.672989 | 0 |
| 7 | 8.635849 | 203.361523 | 13672.091764 | 4.563009 | 303.309771 | 474.607645 | 12.363817 | 62.798309 | 4.401425 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3267 | 8.989900 | 215.047358 | 15921.412018 | 6.297312 | 312.931022 | 390.410231 | 9.899115 | 55.069304 | 4.613843 | 1 |
| 3268 | 6.702547 | 207.321505 | 17246.920347 | 7.708117 | 304.510230 | 329.266007 | 16.217303 | 28.878601 | 3.442983 | 1 |
| 3269 | 11.491011 | 94.812545 | 37188.826022 | 9.263156 | 258.930600 | 439.893618 | 16.172755 | 41.558501 | 4.369264 | 1 |
| 3270 | 6.069616 | 186.659040 | 26138.780191 | 7.747547 | 345.700257 | 415.886955 | 12.067620 | 60.419921 | 3.669712 | 1 |
| 3271 | 4.668102 | 193.681735 | 47580.991603 | 7.166639 | 359.948574 | 526.424171 | 13.894419 | 66.687695 | 4.435821 | 1 |

2011 rows × 10 columns

In [ ]:

jupyter Untitled4 Last Checkpoint: an hour ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

Trusted | Python 3 (ipykernel)

Code

2011 rows × 10 columns

In [27]: `d.isnull()`

Out[27]:

| | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|------|-------|----------|--------|-------------|---------|--------------|----------------|-----------------|-----------|------------|
| 0 | True | False | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | True | False | False | False | False | False |
| 2 | False | False | False | False | True | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False | False |
| ... | | | | | | | | | | |
| 3271 | False | False | False | False | False | False | False | False | False | False |
| 3272 | False | False | False | True | False | False | True | False | False |
| 3273 | False | False | False | False | True | False | False | False | False | False |
| 3274 | False | False | False | False | True | False | False | False | False | False |
| 3275 | False | False | False | False | True | False | False | False | False | False |

3276 rows × 10 columns

In [ ]:

In [28]: d.notnull()

Out[28]:

|  | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | True | True | True | True | True | True | True | True | True |
| 1 | True | True | True | True | False | True | True | True | True | True |
| 2 | True | True | True | True | False | True | True | True | True | True |
| 3 | True | True | True | True | True | True | True | True | True | True |
| 4 | True | True | True | True | True | True | True | True | True | True |
| ... | | | | | | | | | | |
| 3271 | True | True | True | True | True | True | True | True | True | True |
| 3272 | True | True | True | True | False | True | True | False | True | True |
| 3273 | True | True | True | True | False | True | True | True | True | True |
| 3274 | True | True | True | True | False | True | True | True | True | True |
| 3275 | True | True | True | True | False | True | True | True | True | True |

2276 rows × 10 columns

In [ ]:

3275  True  True  True  True  False  True  True  True  True  True

3276 rows × 10 columns

In [29]: `d.fillna(0)`

Out[29]:

|  | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.000000 | 204.890455 | 22791.310501 | 7.300212 | 368.516441 | 564.308654 | 10.379783 | 86.990970 | 2.963135 | 0 |
| 1 | 3.716080 | 129.422921 | 18630.057858 | 6.635246 | 0.000000 | 592.885359 | 15.180013 | 56.329076 | 4.500656 | 0 |
| 2 | 8.099124 | 224.236259 | 19909.541732 | 9.275884 | 0.000000 | 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 | 9.092220 | 181.101509 | 17978.986339 | 6.546600 | 310.135738 | 398.410813 | 11.558279 | 31.997993 | 4.075075 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3271 | 4.668102 | 193.681735 | 47580.991603 | 7.166639 | 359.948574 | 526.424171 | 13.894419 | 66.687695 | 4.435821 | 1 |
| 3272 | 7.808856 | 193.553212 | 17329.802160 | 8.061362 | 0.000000 | 392.449580 | 19.903225 | 0.000000 | 2.798243 | 1 |
| 3273 | 9.419510 | 175.762646 | 33155.578218 | 7.350233 | 0.000000 | 432.044783 | 11.039070 | 69.845400 | 3.298875 | 1 |
| 3274 | 5.126763 | 230.603758 | 11983.809370 | 6.303357 | 0.000000 | 402.883113 | 11.168946 | 77.488213 | 4.708658 | 1 |
| 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509305 | 0.000000 | 327.459760 | 16.140368 | 78.698446 | 2.309149 | 1 |

3276 rows × 10 columns

In [ ]: |

```
Run ▶  ■  C  ⏭  Code
```

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 4 | 9.092223 | 181.101509 | 17978.986339 | 6.546600 | 310.135738 | 398.410813 | 11.558279 | 31.997993 | 4.075075 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3271 | 4.668102 | 193.601735 | 47580.991603 | 7.166639 | 359.948574 | 526.424171 | 13.894419 | 66.687699 | 4.435821 | 1 |
| 3272 | 7.808856 | 193.553212 | 17329.802160 | 8.061362 | 0.000000 | 392.449580 | 19.903225 | 0.000000 | 2.798243 | 1 |
| 3273 | 9.419510 | 175.762646 | 33155.578218 | 7.350233 | 0.000000 | 432.044783 | 11.039070 | 69.845400 | 3.298879 | 1 |
| 3274 | 5.126763 | 230.603758 | 11983.886976 | 6.303357 | 0.000000 | 402.883113 | 11.168945 | 77.488213 | 4.708658 | 1 |
| 3275 | 7.874671 | 195.103239 | 17404.177061 | 7.509005 | 0.000000 | 327.459760 | 16.140368 | 78.698046 | 2.309149 | 1 |

3276 rows × 10 columns

In [33]: d.dtypes

Out[33]:
```
ph                 float64
Hardness           float64
Solids             float64
Chloramines        float64
Sulfate            float64
Conductivity       float64
Organic_carbon     float64
Trihalomethanes    float64
Turbidity          float64
Potability           int64
dtype: object
```

In [ ]:

```
In [34]: d.tail(15)
```

Out[34]:

| | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|---|
| 3261 | 3.620922 | 244.187392 | 24855.033209 | 6.618071 | 356.957873 | 442.070337 | 13.302830 | 59.480204 | 4.754826 | 1 |
| 3262 | 8.372168 | 198.511215 | 28474.292380 | 6.477057 | 319.477187 | 489.868994 | 13.389083 | 35.221200 | 4.624093 | 1 |
| 3263 | 6.920806 | 280.590151 | 24793.525623 | 5.550154 | 332.232177 | 507.723587 | 15.403027 | 51.535057 | 4.013309 | 1 |
| 3264 | 5.893103 | 239.269481 | 20525.666158 | 6.342594 | 341.256362 | 403.617362 | 18.063707 | 63.846319 | 4.390702 | 1 |
| 3265 | 8.197353 | 203.106091 | 27701.794155 | 6.472914 | 328.856838 | 444.012724 | 14.250876 | 62.500205 | 3.351893 | 1 |
| 3266 | 0.372910 | 169.067052 | 14822.745191 | 7.547604 | NaN | 481.525552 | 11.403027 | 30.436151 | 4.906260 | 1 |
| 3267 | 8.989900 | 215.047358 | 15921.412015 | 6.297312 | 312.931022 | 395.410231 | 9.809116 | 55.069304 | 4.613843 | 1 |
| 3268 | 6.702547 | 207.321085 | 17246.920347 | 7.708117 | 304.510280 | 329.266002 | 16.217303 | 28.878921 | 3.442983 | 1 |
| 3269 | 11.491011 | 94.812545 | 37188.026022 | 9.263166 | 258.930600 | 439.893516 | 16.172755 | 41.558521 | 4.369264 | 1 |
| 3270 | 6.069616 | 186.659040 | 26138.780191 | 7.747547 | 345.700257 | 415.886955 | 12.067620 | 60.419971 | 3.669712 | 1 |
| 3271 | 4.668102 | 193.681735 | 47580.991603 | 7.166639 | 359.948574 | 526.424171 | 13.894419 | 66.687695 | 4.435821 | 1 |
| 3272 | 7.808856 | 193.553212 | 17329.802160 | 8.061362 | NaN | 392.449580 | 19.903225 | NaN | 2.798243 | 1 |
| 3273 | 9.419510 | 175.762548 | 33155.578318 | 7.350230 | NaN | 432.044703 | 11.039070 | 69.845400 | 3.298075 | 1 |
| 3274 | 5.126763 | 230.603755 | 11983.869375 | 6.303357 | NaN | 402.883113 | 11.168946 | 77.488213 | 4.708658 | 1 |
| 3275 | 7.874671 | 195.102299 | 17404.177061 | 7.509306 | NaN | 327.459760 | 16.140368 | 78.698446 | 2.309149 | 1 |

```
In [ ]:
```

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 3271 | 4.668102 | 193.601705 | 17500.991603 | 7.165838 | 359.948574 | 526.424171 | 13.894419 | 66.607695 | 4.435021 | 1 |
| 3272 | 7.808656 | 193.553212 | 17329.802165 | 8.061362 | NaN | 392.449580 | 19.903225 | NaN | 2.798243 | 1 |
| 3273 | 9.419510 | 175.762546 | 33155.578215 | 7.350233 | NaN | 432.044783 | 11.039070 | 69.845400 | 3.298375 | 1 |
| 3274 | 5.126763 | 230.603758 | 11983.869170 | 6.303357 | NaN | 402.883112 | 11.168946 | 77.488213 | 4.708658 | 1 |
| 3275 | 7.874671 | 195.102299 | 17404.177261 | 7.509306 | NaN | 327.459760 | 16.140368 | 78.698446 | 2.309149 | 1 |

In [35]: d.head(10)

Out[35]:

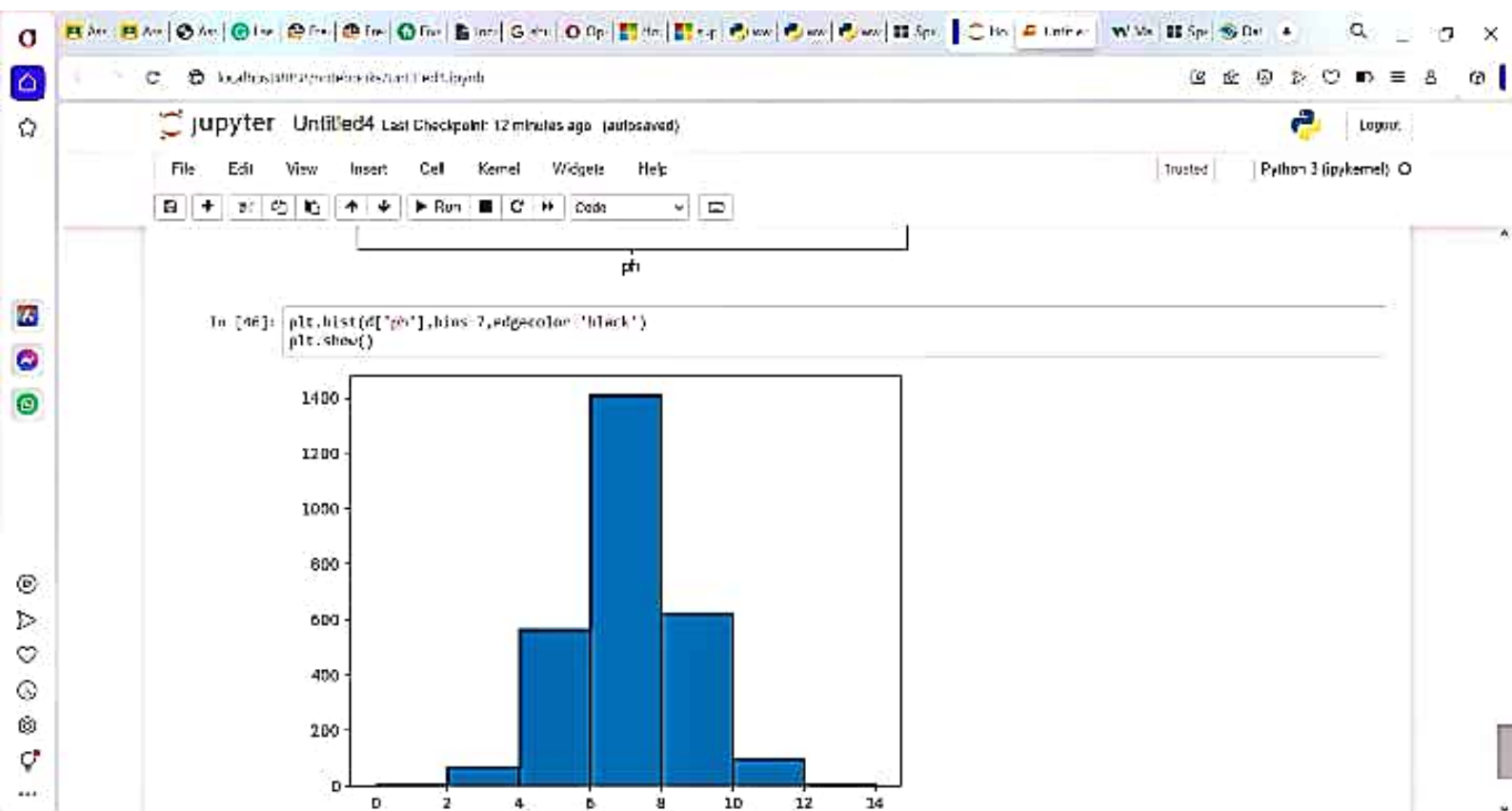| | ph | Hardness | Solids | Chloramines | Sulfate | Conductivity | Organic_carbon | Trihalomethanes | Turbidity | Potability |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NaN | 204.890456 | 20791.318981 | 7.300212 | 368.516441 | 564.308654 | 10.379783 | 86.990970 | 2.963135 | 0 |
| 1 | 3.716080 | 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 |
| 5 | 5.584087 | 188.313324 | 28748.687739 | 7.544869 | 326.678363 | 280.467916 | 8.399735 | 54.917862 | 2.559708 | 0 |
| 6 | 10.223862 | 248.071735 | 28749.716544 | 7.513408 | 391.663396 | 283.651634 | 13.789695 | 84.603556 | 2.672989 | 0 |
| 7 | 8.635849 | 203.361523 | 13672.091764 | 4.563009 | 303.309771 | 474.607645 | 12.363817 | 62.798309 | 4.401425 | 0 |
| 8 | NaN | 118.988579 | 14285.583854 | 7.804174 | 268.646941 | 389.375566 | 12.706049 | 53.928846 | 3.595017 | 0 |
| 9 | 11.180284 | 227.231469 | 25484.508491 | 9.077200 | 404.041635 | 563.885481 | 17.927806 | 71.976601 | 4.370562 | 0 |

In [ ]:

# Performing analysis and Visualization:

- Analyzing and visualizing a dataset is a common and crucial step in data analysis.
  - **Load the data**
  - **Explore the data**
    - Start by looking at the first few rows of your dataset to understand its structure and the type of data it contains.
    - Use functions like **head().**
  - **Data Cleaning**
    - This might involve imputing missing values, removing duplicates, or filtering out extreme outliers.
  - **Descriptive Statistics**
    - Calculate basic statistics like mean, median, standard deviation etc
  - **Data Visualization**
    - Histogram,Scatter plot,Line Charts ,Bar Charts,Box plot,etc

```
import seaborn as sns
```
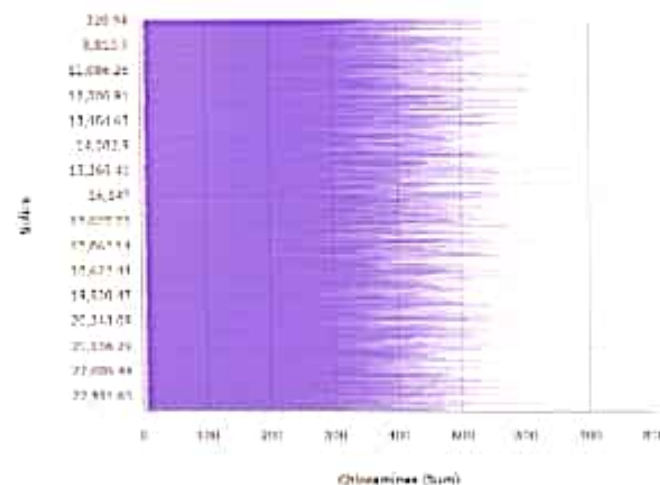
```
In [42]: sns.countplot(x='ph',data=d)
         plt.show()
```

```
In [49]: plt.plot(d['ph'],d['Solids'],marker='o',linestyle='-')
         plt.show()
```

ph

```
In [46]: plt.hist(d['ph'],bins=7,edgecolor='black')
         plt.show()
```

IBM Cognos Analytics          🔵  * New exploration  ⌄                                    30  🔍  💬  ⑦

                                                    Analytics ⤴   Details 🔍   Filters ▽   **Fields** 🔖   Pro

↑  2/2  ⌄   Column 📊   Related 📊  ✐  ⫶   ⟋   ⸘  📊  🔲   Sync ⌄  🗑

Selected sources /

water_potability.csv        +  1

🔍 Search

   📁 Navigation paths         ⁎

⁎  ▦ water_potability.csv

    f⸮ ph

    ▙⸮ Hardness

    ▙⸮ Solids

    f⸮ Chloramines

    ▙⸮ Sulfate

    f⸮ Conductivity

    ▙⸮ Organic_carbon

    ▙⸮ Trihalomethanes

    f⸮ Turbidity

    ▙⸮ Potability

**Chart A**

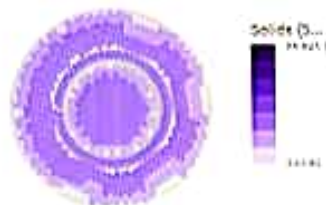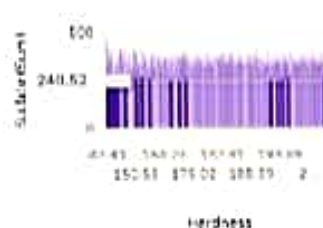**Sulfate and Trihalomethanes hierarchy colored by Solids**  ⓘ

Solids (S...

**Chart B**

**Sulfate by Hardness**  ⓘ

SulfateCount

100

240.52

 

43.43   548.24   557.45   598.96
    150.51  175.02  188.55    2

Hardness

| Summary | Chart A : Solids | Chart B : ⌄ | | Combined |
|---|---|---|---|---|
| Chart perc... data set | 85.43% | 91.73% | | |
| Average | 28,536.04 | 334.31 | ⸱ | |
| Chart total | 61,600,128.89 | 763,894.30 | ⸱ | |
| ⸱ | | | ⸱ | |

**Fields**

📊 Bars

  ⯂ Hardness

    Click or drop data here

# **Length\***                         Required

  ⯂ Sulfate

    Click or drop data here

# y-start

  ⯂ Turbidity

⚙ Target

☰  IBM Cognos Analytics  |  ▦ *New data module  ⌄            ● 🔍 🗩 ⑦ 🔔 👤

Assistant                                                          ↻ ⤢ ✕

▦ water_potability.csv ⓘ  Change

📊 📊 📊 ▦ ❚❚ 🌙 ✈ 🏠 ▥ 📈

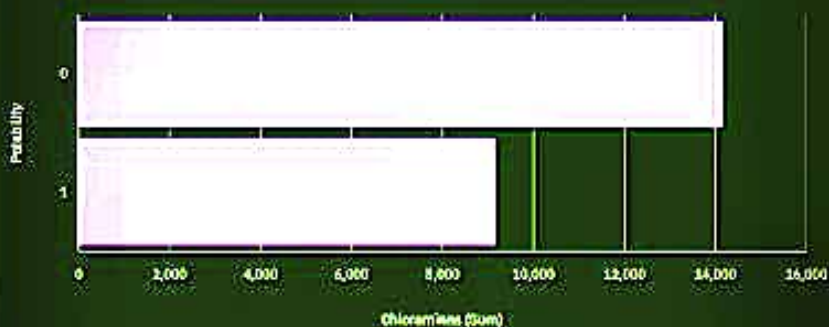### Chloramines by Potability                                    ⓘ        ✴ Top Insights                                              ⓘ



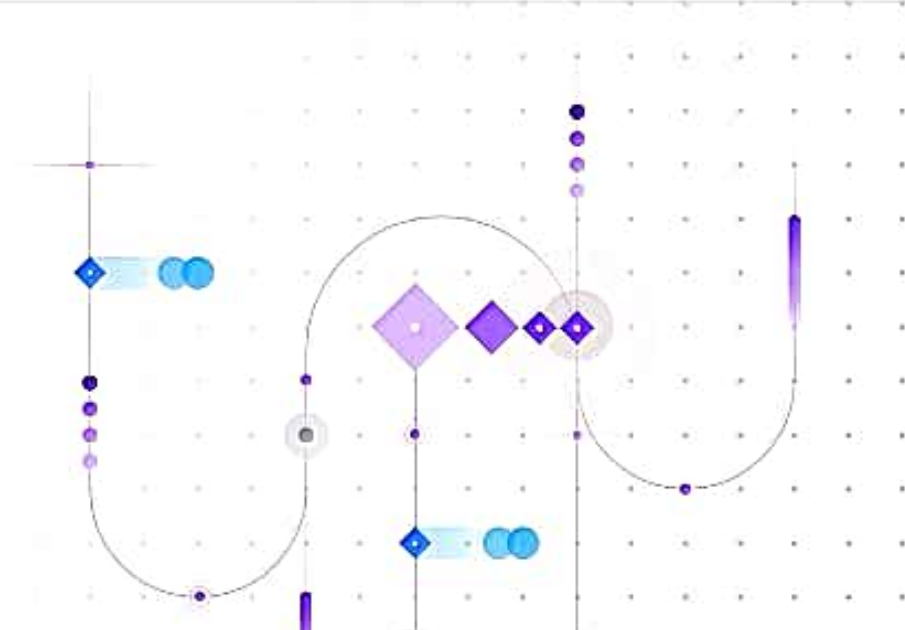|                                                    | 0 exceeds 1 in Chloramines by 5008. |
|                                                    | **Potability 1 has the lowest total Chloramines at over nine thousand.** |
|                                                    | **Potability 0 has the highest total Chloramines at over 14 thousand.** |
|                                                    | Across all values of Potability, the sum of Chloramines is over 23 thousand. |
|                                                    | Chloramines ranges from over nine thousand, when Potability is 1, to over fourteen thousand, when Potability is 0. |

📈