# WATER QUALITY ANALYSIS

## PHASE 3: DATA PREPROCESSING & VISUALIZATION

## INTRODUCTION:

The diversity of research on the water quality monitoring system discussed by Wang, Wang (2009) and Jiang & All (2009). According to Wang, Wang (2009), IoT devices typically in this application usually include large geographical areas and also mobile. The IoT diversification leads to wireless sensing that will revolutionize the river area. The developed system allows for remote monitoring and real-time monitoring of water quality parameters and allowing current status observations and water quality history. Jiang & All (2009) perform a water monitoring system that includes pH analysis, conductivity, dissolved oxygen and temperature levels will be implemented. Alarm sounds will be triggered if there is water pollution or water quality changes. The parameters are measured with landing sensors and data transmitted to base stations via GPRS.

## DATA PREPROCESSING:

Data preprocessing is a crucial step in data analysis and machine learning. It involves cleaning and transforming raw data into a format that is suitable for analysis or model training. The process typically includes:

1.**Data Cleaning:** Removing or handling missing values, correcting errors, and dealing with outliers to ensure data quality.

2.**Data Transformation:** Converting data into a suitable format, such as encoding categorical variables, scaling numerical features, and normalizing data.

3.**Feature Selection:** Identifying and selecting the most relevant features (variables) for analysis or model training to improve efficiency and reduce noise.

4.**Data Reduction:** Reducing the dimensionality of data when dealing with high-dimensional datasets to simplify analysis.

5.**Data Integration**: Combining data from multiple sources or datasets to create a unified dataset for analysis or modeling.

# MISSING VALUES:

Missing values can be handled using various libraries and techniques. Importing necessary libraries called pandas and matplotlib.

## Pandas:

Pandas is an open-source data manipulation and analysis library for Python. It provides data structures and functions that make it easier to work with structured data, such as tables or CSV files. Pandas is particularly useful for data cleaning, transformation, and analysis tasks. It offers two primary data structures: Series (1D labeled arrays) and DataFrame (2D labeled data structures), making it a powerful tool for data processing and analysis.

## Matplotlib:

Matplotlib is a widely-used Python library for creating high-quality static, animated, or interactive visualizations. It provides a wide range of options for creating plots and charts, allowing users to customize every aspect of their visualizations. Matplotlib is versatile and can be used for various types of graphs, including line plots, bar charts, scatter plots, and more. It's often used in combination with other libraries like NumPy and Pandas to visualize data and gain insights from it.



In the above picture missing values are filled using python libraries called pandas and matplotlib.Missing values are filled with "Nan".

# VISUALIZATION:

Creating visualizations can enhance the understanding of complex data. Consider using graphs or charts to represent water quality parameters over time, making trends and fluctuations more apparent. This visual approach facilitates better communication of the analysis results.



# ACCURACY:

Accuracy in water quality analysis is essential to ensure reliable and meaningful results. Here are some key factors to consider for achieving accuracy:

1.**Quality Assurance:** Implement a robust quality assurance/quality control (QA/QC) program that includes
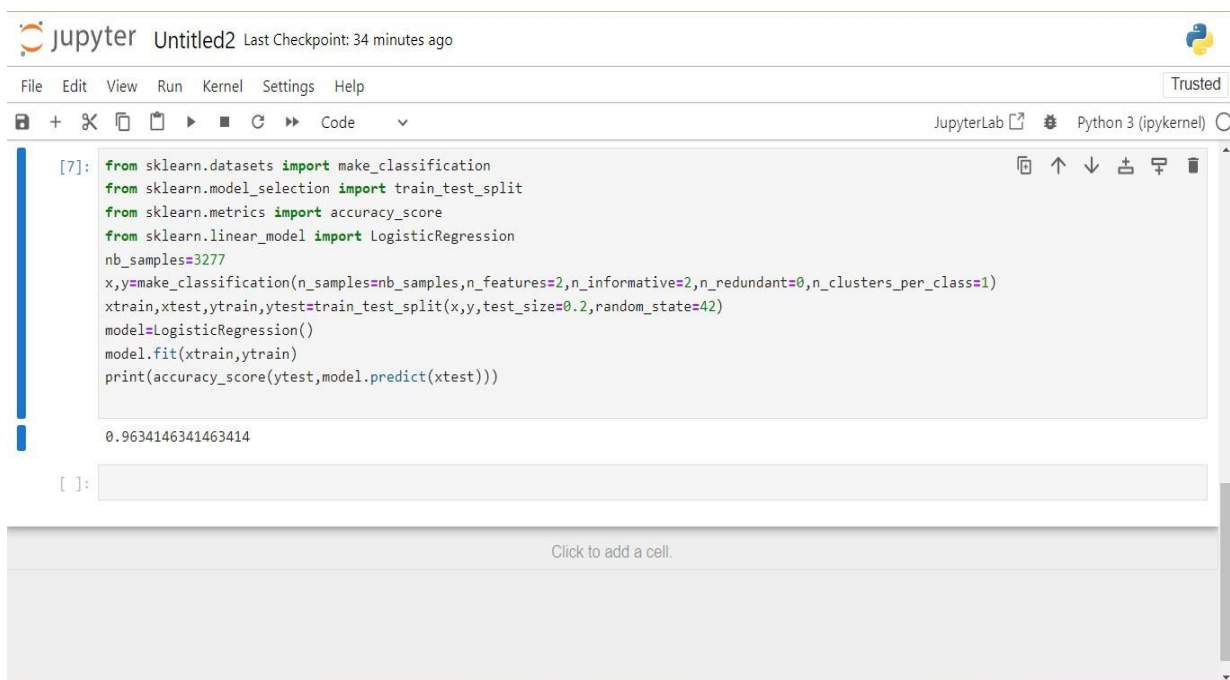
calibration checks, replicates, and certified reference materials.

2.**Proper Sampling:** Ensure that samples are collected using appropriate methods, containers, and at the right locations and times to represent the water source accurately.

3.**Calibration:** Calibrate instruments regularly to maintain accuracy. This includes pH meters, spectrophotometers, and other analytical tools.

4.**Precision:** Consistency in sample handling and analysis methods helps minimize errors and improves accuracy.

5.**Trained Personnel:** Employ trained and experienced personnel to conduct water quality analysis and interpret results correctly.

In the above picture we have calculated the accuracy of our water quality analysis dataset.

## CONSISTENCY:

Consistency in water quality analysis involves using standardized methods and protocols consistently over time. Regular calibration of instruments, proper sample handling, and adherence to established procedures ensure reliability in results. This consistency is essential for tracking changes, identifying trends, and making informed decisions about water management and environmental.



```python
ph_range = (6.5, 8.5)
dissolved_oxygen_range = (5, 10)
temperature_range = (20, 30)
def is_consistent(ph, dissolved_oxygen, temperature):
    if ph_range[0] <= ph <= ph_range[1] and dissolved_oxygen_range[0] <= dissolved_oxygen <= dissolved_oxygen_range[1] and temperature_ra
        return True
    else:
        return False
sample_ph = 7.2
sample_dissolved_oxygen = 8.3
sample_temperature = 25

if is_consistent(sample_ph, sample_dissolved_oxygen, sample_temperature):
    print("The water quality is consistent.")
else:
    print("The water quality is not consistent.")
```

```
The water quality is consistent.
```

## CONCLUSION:

Based on the water quality analysis, it appears that the water meets acceptable standards, ensuring its suitability for consumption and other purposes. Regular monitoring is advisable to maintain these standards.