

# Water quality analysis with python programming

BATCH MEMBER

812121104022:Jeeva Arunkumar

Project title :Water quality analysis

Phase3:Development part 2

Topic:Start building the water quality analysis model  
by python

## Introduction: Water Quality Analysis

- Analysing water quality is one of the key topics of machine learning research.
- In order to train a machine learning model that can determine if a certain water sample is safe or unsafe for eating, we must first understand all the parameters that impact water potability.
- This process is also known as water potability analysis.
- We'll be utilising a Kaggle dataset that includes information on all of the key elements that have an impact on the potability of water for the water quality analysis challenge.
- Before building a model using machine learning to predict whether the water specimen is acceptable or unsafe for eating.
- we must first quickly examine each characteristic of this dataset because all of the elements that determine water quality are crucial.



## **Overview of the process:**

### **1. \*\*Sample Collection\*\*:**

- Select sampling locations that are representative of the water source or system being studied.
- Use clean, non-contaminated containers and equipment for sample collection.
- Collect samples at different depths if applicable (e.g., for lakes and reservoirs).

### **2. \*\*Sample Preservation\*\*:**

- Preserve the samples as necessary to maintain their integrity until analysis. Common preservation methods include:
  - Refrigeration: For most samples, storing at 4°C is sufficient.
  - Chemical preservatives: Adding specific reagents to prevent microbial growth or chemical changes.

### **3. \*\*Physical Characteristics Analysis\*\*:**

- Measure physical parameters, including temperature, turbidity, and conductivity.
- Conduct a visual inspection for color, odor, and sediment.

### **4. \*\*Chemical Characteristics Analysis\*\*:**

- Determine the chemical composition of the water. Common chemical parameters include:
  - pH (acidity or alkalinity)

- Dissolved oxygen (DO)
- Chemical oxygen demand (COD)
- Biochemical oxygen demand (BOD)
- Nutrient concentrations (nitrate, phosphate, etc.)
- Metals (e.g., iron, lead, copper)
- Organic compounds (e.g., pesticides, volatile organic compounds)

## **5. \*\*Biological Characteristics Analysis\*\*:**

- Assess the presence of microorganisms, such as coliform bacteria or fecal coliforms, as indicators of microbial contamination.
- Evaluate the diversity and abundance of aquatic organisms in the water, including algae, plankton, and macroinvertebrates.

## **6. \*\*Toxicological Analysis\*\*:**

- Determine the presence of toxic substances or pollutants using specialized tests for specific contaminants, such as heavy metals, pesticides, or organic pollutants.

## **7. \*\*Data Interpretation\*\*:**

- Compare the obtained data with established water quality standards, guidelines, or regulatory limits to assess the water's suitability for its intended use.
- Identify potential water quality issues or areas of concern.

## **8. \*\*Reporting and Communication\*\*:**

- Compile the results into a comprehensive report.
- Communicate the findings to relevant authorities, stakeholders, and the public as necessary.

#### **9. \*\*Quality Control\*\*:**

- Implement quality control measures to ensure the accuracy and reliability of the data, including the use of certified reference materials and duplicate samples.

#### **10. \*\*Ongoing Monitoring\*\*:**

- Continuously monitor water quality over time to detect trends and assess the effectiveness of any remediation efforts.

#### **11. \*\*Remediation and Management\*\*:**

- Implement measures to improve water quality if issues are identified, such as treating water or addressing pollution sources.

Water quality analysis is an ongoing process that helps protect human health, ecosystems, and water resources. The specific parameters and methods used in analysis may vary depending on the purpose, regulatory requirements, and the characteristics of the water source being studied.

## **Feature engineering for water quality analysis:**

### **1. \*\*Temporal Features\*\*:**

- Time-based features can help capture seasonality and trends. Examples include:

- Day of the week, month, or year.
- Time of day.
- Moving averages or rolling statistics to smooth out data.

### **2. \*\*Lagged Features\*\*:**

- Create lagged features by shifting measurements in time. This can capture delayed effects and autocorrelation. For instance, you might use data from the past few days or weeks to predict the current water quality.

### **3. \*\*Statistical Aggregates\*\*:**

- Calculate statistical aggregates such as mean, median, standard deviation, and percentiles for water quality parameters over different time intervals (e.g., hourly, daily, monthly). These can provide a summary of the data distribution.

#### **4. \*\*Interactions\*\*:**

- Create interaction features by combining two or more water quality parameters. For instance, you might calculate the ratio of nitrate to phosphate concentrations to assess nutrient balance.

#### **5. \*\*Derivative Features\*\*:**

- Compute the rate of change for specific water quality parameters over time. Derivatives can be useful in understanding trends and identifying sudden changes or anomalies.

#### **6. \*\*Geospatial Features\*\*:**

- If you have data from different locations, consider including geospatial features, such as latitude and longitude, distance to pollution sources, or proximity to specific geographical features.

#### **7. \*\*Weather Data Integration\*\*:**

- Incorporate weather-related features (e.g., precipitation, temperature, humidity) as they can

influence water quality. Combining water quality data with weather data can help in understanding how environmental factors impact water parameters.

## **8. \*\*Seasonal Decomposition\*\*:**

- Decompose the time series data into trend, seasonality, and residual components. These components can be used as separate features for analysis.

## **9. \*\*Time Since Last Event\*\*:**

- Calculate the time since significant water quality events or occurrences (e.g., heavy rainfall, pollution incidents) to assess the persistence of impacts.

## **10. \*\*Historical Averages and Trends\*\*:**

- Compute historical averages and trends for each water quality parameter, providing context for the current measurement.

## **11. \*\*Frequency Domain Analysis\*\*:**



- Use techniques such as Fourier transforms to extract frequency domain features that may reveal periodic patterns in the data.

## **12. \*\*Principal Component Analysis (PCA)\*\*:**

- Reduce the dimensionality of the data by using PCA to create new features that capture the most significant variations in water quality parameters.

## **13. \*\*Cross-Correlation\*\*:**

- Calculate cross-correlation between different water quality parameters to identify relationships and dependencies between them.

## **14. \*\*Categorical Encoding\*\*:**

- Encode categorical variables like water source type, treatment methods, or regulatory zones into numerical features using techniques like one-hot encoding.

## **15. \*\*Target Transformation\*\*:**

- Apply transformations to the target variable, such as log transformation, to make it more amenable to modeling and account for skewed distributions.

## **16. \*\*Feature Scaling and Normalization\*\*:**

- Normalize or scale features to ensure they are on a common scale, especially when using machine learning algorithms sensitive to feature scales.

The choice of feature engineering techniques depends on the specific characteristics of your water quality dataset and the goals of your analysis. Experiment with different features and assess their impact on model performance to identify the most informative features for your water quality analysis.

Model evaluation:

### **1. \*\*Data Splitting\*\*:**

- Split the dataset into training, validation, and test sets to train the model, tune hyperparameters, and evaluate its performance on unseen data.

### **2. \*\*Cross-Validation\*\*:**

- Use techniques like k-fold cross-validation to assess the model's robustness and generalization ability by repeatedly splitting the data into training and validation subsets.

### 3. **\*\*Performance Metrics\*\***:

#### a. **\*\*Regression Metrics\*\***:

- Mean Absolute Error (MAE): Measures the average absolute difference between predicted and actual values.

- Root Mean Square Error (RMSE): Provides a measure of the model's prediction error by taking the square root of the mean of squared differences.

- R-squared ( $R^2$ ) or Coefficient of Determination: Indicates the proportion of variance explained by the model.

#### b. **\*\*Classification Metrics\*\***:

- Accuracy: Measures the proportion of correctly classified instances.

- Precision: Calculates the ratio of true positive predictions to the total positive predictions.

- Recall (Sensitivity or True Positive Rate): Measures the ability of the model to identify positive instances.

- F1-Score: The harmonic mean of precision and recall, balancing precision and recall.

- Area Under the Receiver Operating Characteristic Curve (AUC-ROC): Evaluates the model's ability to distinguish between classes in binary classification problems.

#### **4. \*\*Confusion Matrix\*\*:**

- For classification tasks, create a confusion matrix to visualize the true positives, true negatives, false positives, and false negatives. This helps in understanding model performance, especially when dealing with imbalanced datasets.

#### **5. \*\*Residual Analysis\*\*:**

- For regression tasks, analyze the model residuals (the differences between predicted and actual values) to check for patterns or biases. A well-behaved residual plot should show random scattering around zero.

#### **6. \*\*Feature Importance\*\*:**

- Assess the importance of each feature in the model using techniques like feature importance scores or permutation importance to understand which variables are most influential in water quality predictions.

## 7. **\*\*Model Selection\*\***:

- Compare the performance of different models (e.g., linear regression, decision trees, random forests, neural networks) to identify the most suitable model for the specific water quality analysis task.

## 8. **\*\*Overfitting and Underfitting Analysis\*\***:

- Check for signs of overfitting (the model is too complex and fits noise) or underfitting (the model is too simple to capture the underlying patterns) by examining training and validation performance.

## 9. **\*\*Bias and Fairness Evaluation\*\***:

- Assess the model for potential biases or unfairness, especially when water quality analysis has implications for environmental justice and social equity.

## 10. **\*\*External Validation\*\***:

- If available, compare model predictions with independent data sources or conduct field tests to validate the model's performance in real-world conditions.

## **11. \*\*Error Analysis\*\*:**

- Analyze the specific types of errors the model makes and investigate potential causes. This can guide improvements in data quality and model design.

## **12. \*\*Model Interpretability\*\*:**

- Ensure that the model is interpretable, especially if the results need to be communicated to stakeholders. Techniques like SHAP values and feature importance can help in understanding model predictions.

## **13. \*\*Cost-Benefit Analysis\*\*:**

- Consider the costs associated with model errors and the benefits of accurate predictions. This can help in selecting an appropriate model threshold or decision boundary.

The choice of evaluation techniques and metrics should be tailored to the specific objectives of the water quality analysis and the type of models being used. Regular model evaluation is essential to ensure that the model

remains accurate and reliable over time, especially in dynamic environmental conditions.

## Visualization:

Model evaluation is a crucial step in water quality analysis to assess the performance of predictive models and determine their reliability in estimating or classifying water quality parameters. Below are some common techniques and metrics for evaluating models in the context of water quality analysis:

### 1. **Data Splitting**:

- Split the dataset into training, validation, and test sets to train the model, tune hyperparameters, and evaluate its performance on unseen data.

### 2. **Cross-Validation**:

- Use techniques like k-fold cross-validation to assess the model's robustness and generalization ability by repeatedly splitting the data into training and validation subsets.

### 3. **Performance Metrics**:

#### a. **Regression Metrics**:



- Mean Absolute Error (MAE): Measures the average absolute difference between predicted and actual values.

- Root Mean Square Error (RMSE): Provides a measure of the model's prediction error by taking the square root of the mean of squared differences.

- R-squared ( $R^2$ ) or Coefficient of Determination: Indicates the proportion of variance explained by the model.

b. **\*\*Classification Metrics\*\***:

- Accuracy: Measures the proportion of correctly classified instances.

- Precision: Calculates the ratio of true positive predictions to the total positive predictions.

- Recall (Sensitivity or True Positive Rate): Measures the ability of the model to identify positive instances.

- F1-Score: The harmonic mean of precision and recall, balancing precision and recall.

- Area Under the Receiver Operating Characteristic Curve (AUC-ROC): Evaluates the model's ability to distinguish between classes in binary classification problems.

#### 4. **\*\*Confusion Matrix\*\***:

- For classification tasks, create a confusion matrix to visualize the true positives, true negatives, false positives, and false negatives. This helps in understanding model performance, especially when dealing with imbalanced datasets.

#### 5. **\*\*Residual Analysis\*\***:

- For regression tasks, analyze the model residuals (the differences between predicted and actual values) to check for patterns or biases. A well-behaved residual plot should show random scattering around zero.

#### 6. **\*\*Feature Importance\*\***:

- Assess the importance of each feature in the model using techniques like feature importance scores or permutation importance to understand which variables are most influential in water quality predictions.

#### 7. **\*\*Model Selection\*\***:

- Compare the performance of different models (e.g., linear regression, decision trees, random forests, neural networks) to identify the most suitable model for the specific water quality analysis task.

#### 8. **\*\*Overfitting and Underfitting Analysis\*\***:

- Check for signs of overfitting (the model is too complex and fits noise) or underfitting (the model is too simple to capture the underlying patterns) by examining training and validation performance.

#### 9. **\*\*Bias and Fairness Evaluation\*\***:

- Assess the model for potential biases or unfairness, especially when water quality analysis has implications for environmental justice and social equity.

#### 10. **\*\*External Validation\*\***:

- If available, compare model predictions with independent data sources or conduct field tests to validate the model's performance in real-world conditions.

#### 11. **\*\*Error Analysis\*\***:

- Analyze the specific types of errors the model makes and investigate potential causes. This can guide improvements in data quality and model design.

## 12. **\*\*Model Interpretability\*\***:

- Ensure that the model is interpretable, especially if the results need to be communicated to stakeholders.

Techniques like SHAP values and feature importance can help in understanding model predictions.

## 13. **\*\*Cost-Benefit Analysis\*\***:

- Consider the costs associated with model errors and the benefits of accurate predictions. This can help in selecting an appropriate model threshold or decision boundary.

## PROGRAM:

```
import plotly.graph_objs as go

index_vals =

data['Potability'].astype('category').cat.codes

fig = go.Figure(data=go.Splom(
    dimensions=[dict(label='ph',
    values=data['ph']),
    dict(label='Hardness',
    values=data['Hardness']),
    dict(label='Solids',
    values=data['Solids']),
    dict(label='Chloramines',
    values=data['Chloramines'])
    dict(label='Sulfate',
    values=data['Sulfate']),
    dict(label='Conductivity',
    values=data['Conductivity']),
    dict(label='Organic_carbon',
    values=data['Organic_carbon']),
    dict(label='Trihalomethanes',
```

```
values=data['Trihalomethanes']),  
dict(label='Turbidity',values=data['Turbidity'])),  
showupperhalf=False,  
text=data['Potability'],  
marker=dict(color=index_vals,  
showscale=False,  
line_color='white', line_width=0.5)  
))  
fig.update_layout(  
title='Water Quality',  
width=1000,  
height=1000,  
)  
fig.show().
```

## CONCLUSION:

In conclusion, water quality analysis is a vital component of environmental monitoring, public health, and the sustainable management of water resources. It encompasses a comprehensive process of assessing the physical, chemical, and biological characteristics of water to determine its suitability for various purposes.



