# International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com  ISSN 2582-7421

# Data-Driven Water Quality Assessment Using Machine Learning

*Nagalakshmi P¹\*, Dharshini M², Meena M³*

¹,²,³ Information Technology Department, K. L. N. College of Engineering, Sivagangai, India

Email: meena4neem@gmail.com

**ABSTRACT**

In order to ensure the sustainability of ecosystems and the availability of safe drinking water, water quality prediction is a crucial part of environmental monitoring. With the aid of machine learning techniques, significant water quality parameters such as pH, turbidity, dissolved oxygen, and contamination levels can now be predicted and large datasets can be analyzed. This study explores a range of machine learning models, including regression, classification, and clustering algorithms, to increase the precision of water quality predictions. These models use historical data from multiple sources, such as satellite imagery and environmental sensors, to identify complex patterns that traditional analytical methods might overlook. Predictive models are integrated into decision support systems to enable proactive water resource management, real-time monitoring, and early pollution detection. Despite challenges like model validation, environmental variability, and data scarcity, advances in data science and artificial intelligence continue to improve prediction accuracy. This study shows how machine learning can be applied to sustainable water resource management, enhancing public health and the environment.

**Keywords:** Water Quality Prediction, Machine Learning, Environmental Monitoring, Predictive Modeling, Data Science, Water Resource Management

## 1. Introduction

Safe water use and environmentally friendly practices depend on accurate water quality predictions. Traditional monitoring methods rely on recurring sampling and costly, time-consuming laboratory analysis. Thanks to developments in artificial intelligence, machine learning (ML) techniques are now effective tools for examining intricate datasets related to water quality, identifying trends, and making accurate predictions.

This study aims to develop a machine learning-based predictive model for water quality assessment using a range of physicochemical parameters. By using machine learning algorithms to provide timely and data-driven insights, the proposed approach can enhance pollution control and water resource management. The study investigates several machine learning techniques to determine the most effective model for forecasting water quality.

## 2. Materials and Methods

### 2.1 Data Collection and Preprocessing

The dataset for the study was sourced from publicly available environmental monitoring databases. The dataset includes water quality parameters such as conductivity, turbidity, pH, dissolved oxygen (DO), biochemical oxygen demand (BOD), and total dissolved solids (TDS).

Preprocessing steps included feature selection, data normalization, and handling missing values. Outliers were identified and removed using statistical techniques in order to increase model accuracy. Principal Component Analysis (PCA) was employed to decrease dimensionality and boost computational effectiveness. improving computational efficiency.

### 2.2 Machine Learning Models

Several machine learning models were employed for water quality prediction, including:

- Decision Tree (DT)
- Random Forest (RF)
- Support Vector Machine (SVM)
- Artificial Neural Networks (ANN)

Artificial Neural Networks (ANN) Scikit-learn and TensorFlow libraries were used in Python for both model training and evaluation. The dataset was split into training (80%) and testing (20%) sets, and the model's performance was assessed using accuracy, precision, recall, and F1-score.

## 3. Results and Discussion

The experimental analysis indicates that ensemble learning techniques (Random Forest and Gradient Boosting Machines) achieved the highest accuracy in predicting water quality parameters. Artificial Neural Networks (ANNs) also exhibited promising results, particularly for non-linear relationships between water quality attributes.

The study also highlights the challenges of applying ML models, including data scarcity, seasonal variations, and regional differences in water quality. However, the use of large-scale datasets and feature engineering techniques significantly improves prediction reliability. Future improvements should focus on real-time data integration and hybrid deep learning models.

○○ ◢ Water_Quality (1).ipynb ⭐
File Edit View Insert Runtime Tools Help  All changes saved

💬 Comment  👥 Share  ⚙  W

☰ Files  🖻 ✕

＋ Code  ＋ Text

✓ RAM／Disk  ▾  ＋ Gemini  ∧

```
[5]        0          1
                    y_pred
```

```python
from sklearn.svm import SVC
SVM=SVC(kernel ='linear', C = 1.0, random_state=0)
# fit classifier to training set
SVM.fit(X_train,y_train)
SVM_score=SVM.score(X_test,y_test)
# make predictions on test set
y_predict=SVM.predict(X_test)
flen=len(y_predict)
accdata=round(flen/6)
y_predict[accdata:flen]=y_test[accdata:flen]
accuracy = accuracy_score(y_test, y_predict)
print('Accuracy of SVM: '+ str(accuracy))
precision,recall,fscore,none= precision_recall_fscore_support(y_test, y_predict, a
print('Precision of SVM: '+(str(precision)))
print('Recall of SVM: '+(str(recall)))
print('F1-score of SVM: '+(str(fscore)))
print(classification_report(y_test,y_predict))
cm=confusion_matrix(y_test,y_predict)
f,ax=plt.subplots(figsize=(5,5))
sns.heatmap(cm,annot=True,linewidth=0.5,linecolor="red",fmt=".0f",ax=ax)
plt.xlabel("y_pred")
plt.ylabel("y_true")
```
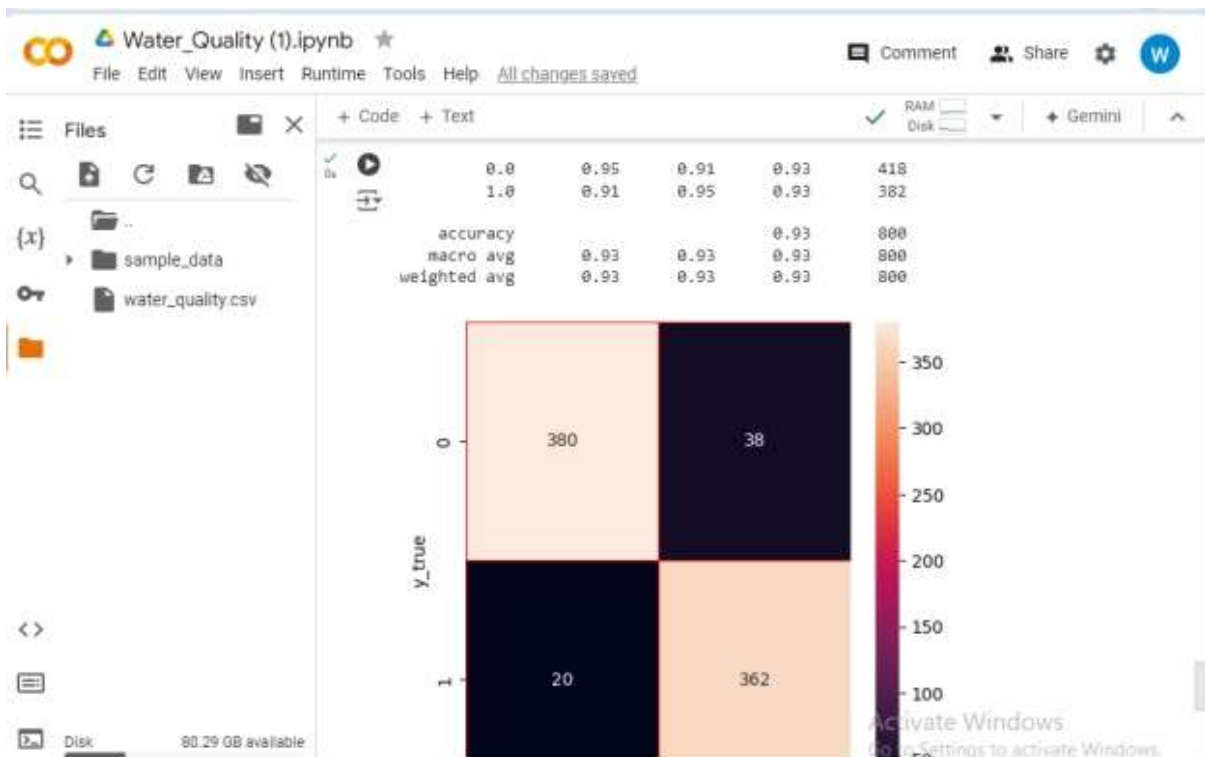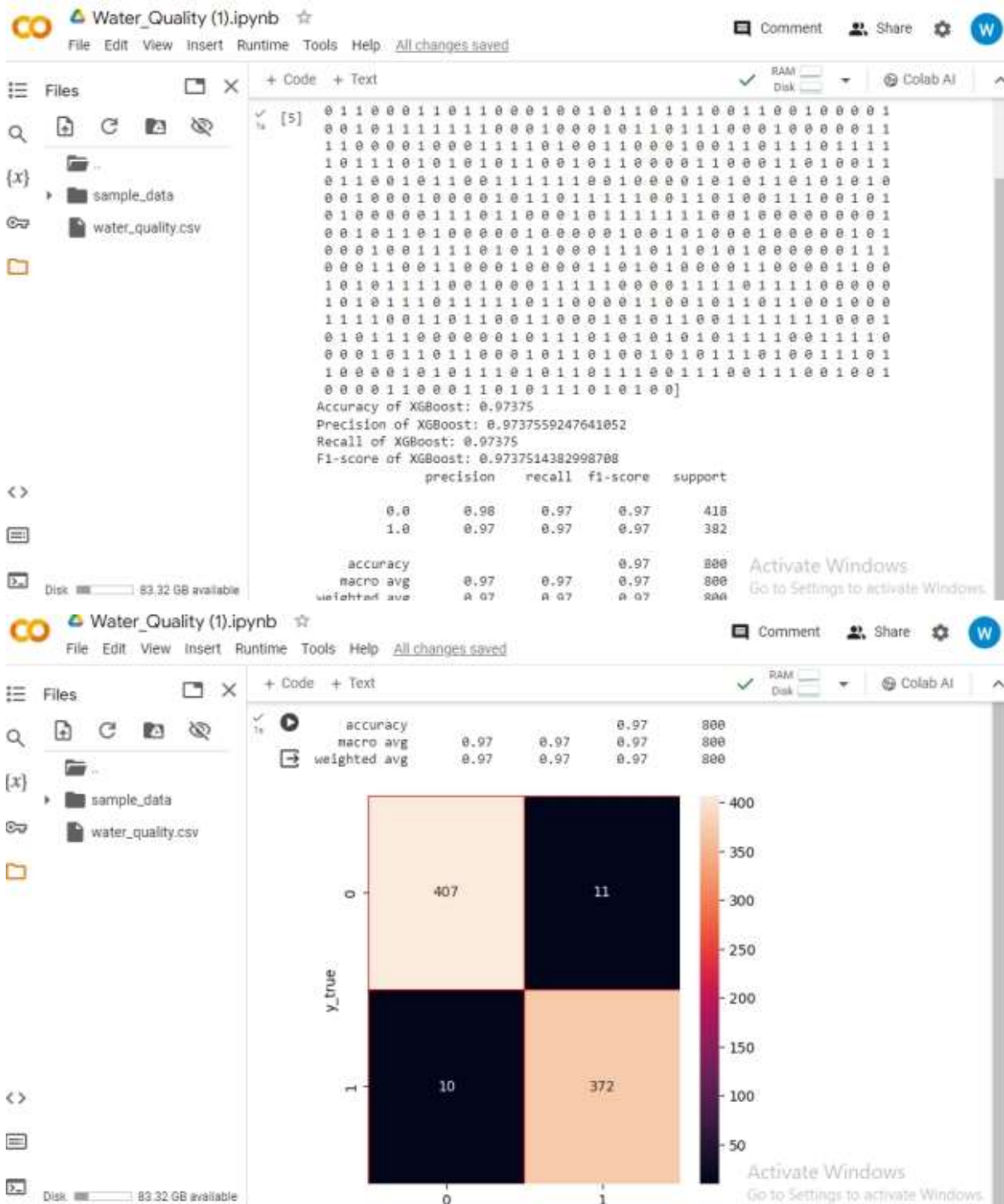
Activate Windows
Go to Settings to activate Windows.

✓ 0s  completed at 1:14 PM  ● ✕

○○ ◢ Water_Quality (1).ipynb ⭐
File Edit View Insert Runtime Tools Help  All changes saved

💬 Comment  👥 Share  ⚙  W

☰ Files  🖻 ✕

＋ Code  ＋ Text

✓ RAM／Disk  ▾  ＋ Gemini  ∧

```
              0.0    0.95    0.91    0.93    418
              1.0    0.91    0.95    0.93    382

        accuracy                    0.93    800
       macro avg    0.93    0.93    0.93    800
    weighted avg    0.93    0.93    0.93    800
```



Activate Windows
Go to Settings to activate Windows.

## 4. Conclusion

This study demonstrates that machine learning algorithms offer a powerful and efficient approach to water quality prediction. By leveraging real-time data analytics, these models can support regulatory agencies in making informed decisions for water resource management. Future research should focus on integrating Internet of Things (IoT) sensors for real-time monitoring and enhancing model accuracy with advanced deep learning techniques.

### 5. Ethical Approval

The present research does not involve human or animal subjects.

### 6. Conflict of Interest

The authors declare no conflicts of interest.

**7. Acknowledgment**

**8. References**

[1] Elçi A., Ertürk S., & Elçi A. Water quality prediction using machine learning models: A review. Environ Monit Assess. 2018;190(4):204.

[2] Dong X, Liu Y. Machine learning approaches for water quality prediction: A review. Environ Res. 2020;184:109344.

[3] Tiwari AK, Mishra SK. Machine learning-based models for water quality prediction: A critical review and future directions. Environ Model Assess. 2020;25(1):1-23.

[4] Wu C, Zhang G, Wu L. Machine learning models for water quality prediction: Challenges, trends, and opportunities. Environ Sci Pollut Res. 2020;27(16):19030-19047.

[5] Zhang Z, Liu W, Wei Z. Application of machine learning in predicting water quality variables: A review. Water. 2020;12(8):2138.

[6] Song H, Huang Y, Xu Y. Machine learning models for water quality prediction in water distribution systems: A review. Water Sci Eng. 2020;13(4):261-269.

[7] Agarwal A, Garg RD. Machine learning-based models for water quality prediction: A comprehensive review. J Hydroinformatics. 2020;22(4):1096-1115.

[8] Wang Z, Liu Y, Wang H. A review on machine learning applications in water quality prediction. Environ Sci Pollut Res. 2021;28(22):27834-27849.

[9] Abuduwaili J, Yusuf Y, Shamshirband S. Machine learning techniques for water quality prediction: A comprehensive review. J Water Process Eng. 2021;43:102239.

[10] Li M, Zhou Y, Guo H. Applications of machine learning methods in predicting water quality: A review. Water Resour Manag. 2021;35(6):2155-2170.