

# **CONFERENCE PROCEEDINGS, JANUARY 2025** Published by the Faculty of Science (FSC), FULafia

Print ISSN: 2354–3388 Online ISSN: 2315–7275 DOI: https://doi.org/10.62050/fscp2024.568

# Machine Learning for Detecting Microbial and Chemical Contaminants in Sachet Water

Rahma Umar Tafida<sup>⊠</sup>, Umar Baba Umar, Hamzat O Aliyu & S. A. Adepoju

Department of Computer Science, Federal University of Technology, Minna, Nigeria

⊠rahma.umar@st.futminna.edu.ng

**bstract:** This study addresses the critical challenge of detecting microbial and chemical contaminants in sachet water in Nigeria using machine learning (ML) techniques. Traditional methods for water quality assessment are often time-consuming, costly, and ill-suited for real-time monitoring, particularly in resource-limited settings. We propose a novel approach that leverages supervised ML algorithms, including Gradient Boosting (GBC) and Random Forest (RF), to predict water potability based on an augmented dataset of 20 parameters, encompassing both microbial contaminants (e.g., *Escherichia coli, Salmonella*) and chemical contaminants (e.g., lead, arsenic). The dataset was enhanced using synthetic data generation techniques to address gaps in the original dataset, which lacked comprehensive coverage of critical contaminants. Our results demonstrate that the Gradient Boosting Classifier (GBC) achieves an accuracy of **99.8%** and an F1 score of **99.7%** on the augmented dataset, significantly outperforming other models. Feature importance analysis revealed that *Escherichia coli, Salmonella*, and lead were the most critical predictors of water potability, aligning with public health concerns. This study highlights the potential of ML for enhancing water quality monitoring, offering a scalable and cost-effective solution to mitigate waterborne diseases in regions like Nigeria, Nigeria. Future work will focus on integrating real-time sensor data and validating the model in real-world scenarios to further improve its applicability and impact.

Keywords: Sachet water, machine learning, microbial contaminants, chemical contaminants

# ntroduction

In Nigeria, the sachet water industry has seen remarkable growth due to persistent water scarcity and concerns about the safety of municipal tap water. Commonly referred to as "pure water," sachet water is often perceived as a safer alternative [1]. However, recent research has revealed consistent microbial and chemical contamination within these products, raising significant concerns about their true safety [2]. Waterborne diseases such as cholera, typhoid, and diarrhea remain critical public health challenges in Nigeria, largely driven by contaminated water sources [3]. Studies have identified harmful contaminants including Escherichia coli, Salmonella, cholerae, heavy metals, and residual Vibrio disinfectants in sachet water across the country [4, 5]; underscoring the urgent need for rigorous water quality monitoring and real-time detection tools [6].

According to the World Health Organization (WHO), poor sanitation, including the widespread practice of open defecation and improper disposal of human waste, significantly contributes to the spread of waterborne diseases [7, 8]. UNICEF data further highlights that water-related diseases are leading causes of mortality among children under five in Nigeria, where only 56% of the population has access to safe drinking water and a mere 13% have access to basic sanitation services. Over 71% of the population practices open defecation, exacerbating health risks UNICEF [9].

Traditional water quality monitoring methods such as laboratory-based chemical analyses, microbiological testing, and physical inspections are time-consuming, labor-intensive, and require specialized expertise [10]. Microbial testing through culture methods can take several days, delaying necessary interventions [11, 12]; while chemical analyses using titration and chromatography are often unsuitable for rapid or realtime monitoring [13]. These limitations hinder timely detection of water quality issues, thereby increasing the public health risks associated with sachet water consumption [14].

Recent advances in machine learning and stochastic mathematics offer promising alternatives for enhancing the accuracy and scalability of water quality predictions [15]. Emerging technologies like remote sensing and wireless sensor networks have contributed to real-time water monitoring systems, although much of the existing research remains focused on periodic rather than continuous monitoring [16]. The Water Quality Index (WQI), a widely used metric, simplifies complex water quality data into a single, interpretable score, supporting effective decision-making [17]. Machine learning models have the potential to improve WQI predictions, even when traditional methods struggle with incomplete data or delayed results.

The increasing consumption of sachet water in Nigeria, spurred by concerns over water scarcity and contamination, has led to significant uncertainties regarding the actual purity and safety of these products [18]. Waterborne diseases, primarily caused by microbial and chemical contaminants, pose severe public health risks, contributing to high morbidity and mortality rates, particularly among vulnerable populations [19]. The critical research problem is the lack of effective real-time monitoring tools specifically designed for detecting these contaminants in sachet water. Traditional assessment methods are often inadequate, lacking the necessary sensitivity and rapid response needed for effective water quality management in rapidly changing environments [20].

In Nigeria, where waterborne diseases such as cholera and typhoid are prevalent and represent significant public health challenges, the urgency for innovative solutions to monitor water quality becomes even more pronounced. Existing methods for water quality assessment, while reliable, are typically labor-intensive and time-consuming, hindering timely responses to contamination events. For example, microbial testing through culture methods can take several days to yield results, complicating efforts to address immediate health risks [21]. Additionally, chemical analysis techniques, such as titration and chromatography, are not conducive to real-time monitoring, further exacerbating the issue.

Many researchers have explored machine learning (ML) models for water quality prediction, aiming to enhance accuracy and provide valuable insights into water potability. Haghiabi *et al.* [22] assessed various artificial intelligence techniques, including artificial neural networks (ANN), group data processing methods (GMDM), and support vector machines (SVM), to predict water quality parameters in the Tireh River, Iran. Their findings indicated that both ANN and SVM models performed well in predicting water quality constituents, with SVM demonstrating the highest accuracy.

In a study on the Kelantan River using data from 2005 to 2020, a set of ML models was investigated to predict water quality classification. Using 13 physical and chemical water quality parameters, Ahmed *et al.* evaluated seven ML models: SVM, ANN, Decision Tree (DT), K-Nearest Neighbors (KNN), Naive Bayes (NB), Random Forest (RF), and Gradient Boosting (GB). Among these, the ensemble model with Gradient Boosting showed superior performance, achieving an accuracy of 94.9%, a sensitivity of 80.0%, and an F-measure of 86.49%, while minimizing classification error [23].

El Bilali and Taleb developed ML models to predict irrigation water quality in semi-arid areas using conductivity and pH as input parameters [24]. Their study demonstrated the models' potential for application in areas where these parameters are critical for irrigation purposes. Similarly, Aldhyani *et al.* explored water quality classification using SVM, KNN, and Naive Bayes, with results indicating that these models effectively predict water quality based on seven important water quality parameters [25].

Another study by Lu and Ma focused on hybrid ML models to improve short-term water quality predictions. Using data from Nainital Lake, they applied extreme gradient boosting and random forest algorithms, showing that Random Forest was the most efficient for regression tasks. However, for classification tasks, Stochastic Gradient Descent, RF, and SVM performed similarly, proving effective in predicting water quality [26].

Dritsas and Trigka evaluated several ML models, such as Naive Bayes, kNN, Logistic Regression (LR), and

tree-based classifiers, using physiochemical and microbiological parameters to classify water as safe or unsafe [27]. They employed SMOTE (Synthetic Minority Oversampling Technique) and used 10-fold cross-validation, with a stacking classification model that achieved high performance (98.1% accuracy, 100% precision, 98.1% recall, and an AUC of 99.9%).

Wang *et al.* investigated SVM, RF, XGBoost, Multilayer Perceptron (MLP), and Long Short-Term Memory (LSTM) for water quality prediction [28]. They found that SVM was highly effective, showing robust generalization capabilities and high prediction accuracy, while MLP was well-suited for nonlinear modeling. However, RF and XGBoost performed less effectively in their study.

Despite these advances, the existing datasets used in these studies often focus on a limited set of water quality indicators. In our analysis, we observed similar limitations in the sachet water dataset, which initially only covered pH, hardness, solids, chloramines, sulfate, conductivity, organic carbon, trihalomethanes, and turbidity. These parameters, while relevant, do not fully represent the contaminants impacting sachet water quality. To address this gap, our research expanded the dataset to include additional parameters that cover both microbial and chemical contaminants. This broader dataset enables more comprehensive ML modeling, improving the detection of water contaminants that directly impact public health, as supported by research on contaminants like coliform bacteria, Escherichia *coli*, lead, arsenic, and others frequently found in water sources.

# aterials and Methods

The data collection process for this study has been completed by conducting a comprehensive literature review. In the course of analyzing the existing water quality dataset, it was observed that the existing dataset, while useful, had significant limitations in terms of the parameters it covered. The original dataset provided values for factors such as pH, hardness, solids, chloramines, sulfate, conductivity, organic carbon, trihalomethanes, and turbidity. However, these parameters do not fully account for all the contaminants that can affect sachet drinking water, as documented in various studies and literature. To address this limitation, we aimed to increase the scope of the dataset by incorporating additional parameters covering both microbial and chemical contaminants that are critical to water quality. This methodology involved identifying and synthesizing existing research on microbial and chemical contaminants found in sachet drinking water in Nigeria. Based on documented health impacts and contaminants found in sachet water, we augmented the dataset with the following new parameters:

**Microbial Contaminants:** Coliform Bacteria, Escherichia coli, Salmonella spp., Vibrio cholerae, Shigella spp., Pseudomonas aeruginosa, Enterococci.

**Chemical Contaminants:** Nitrate  $(NO_2^{-})$ , Lead (Pb),

Arsenic (As), Chlorine Residues (THMs), Fluoride (F<sup>-</sup> ), Aluminum (Al), Zinc (Zn), Copper (Cu), Iron (Fe), Total Dissolved Solids (TDS), and Sulfate (SO  $_{4}^{2-}$ ).

The original Water quality dataset "Water Quality Dataset" was used in the study. The open-access dataset was accessed from the Kaggle website on March 30, 2023. The dataset contains physical and chemical measurements of water quality for 3276 different water bodies. These measurements include nine different property variables: pH value, hardness, sulfate, conductivity, organic solids. carbon, trihalomethanes, turbidity, and potability.

To align the new dataset with the existing dataset's structure, we needed to ensure that the new parameters covered all 3,281 rows of the existing dataset. Since the original dataset did not contain any values for these new parameters, we utilized the synthpop library to augment the dataset by generating synthetic data for the additional columns. The synthpop library is a powerful tool designed for generating synthetic datasets that closely resemble real datasets.

the dataset is likely to be imbalanced, with certain contaminants (e.g., E. coli, Salmonella) appearing more frequently than others (e.g., Vibrio cholerae, heavy metals). This imbalance can lead to biased models that perform well on the majority class (common contaminants) but poorly on the minority class (rare contaminants). For example, a model trained on an imbalanced dataset might achieve high overall accuracy but fail to detect rare but dangerous contaminants, which could have severe public health implications.

To address this issue, SMOTE can be applied to generate synthetic samples for the minority class. SMOTE works by selecting a sample from the minority class and finding its k-nearest neighbors.

The dataset was then partitioned into 80% training and 20% testing sets to ensure a balanced evaluation and improve model reliability.

# esults and Discussion

This section presents the results of our experiments, focusing on the performance of machine learning models in detecting microbial and chemical contaminants in sachet water. The experiments are carried out using the jupyter notebook version (6.4.6). Jupyter notebook makes it easier to run and write Python scripts. We evaluate the models using key metrics such as accuracy, precision, recall, F1

score, and AUC (Area Under the Curve). The results are presented in a clear and coherent manner to demonstrate the effectiveness of our approach. The proposed models performance is compared to that of numerous existing models. The classification models performance was assessed using assessment criteria such as accuracy, recall, precision, F1 score; enhancing classification performance.

#### Experiment Original I: Dataset with **SMOTE+TOMEK**

The original dataset consisted of nine parameters: pH, hardness, solids, chloramines, sulfate, conductivity, organic carbon, trihalomethanes, and turbidity. We the Synthetic Minority applied Over-sampling Technique (SMOTE) combined with Tomek Links to balance the dataset, addressing the issue of class imbalance. The Gradient Boosting Classifier (GBC) was trained on this dataset, and its performance is summarized in Table 1.

#### **Table 1: Classification report**

Parameter	Precision	Recall	F1-Score	Support
0	0.79	0.79	0.79	541
1	0.80	0.78	0.95	512
Accuracy			0.97	1053
Macro Avg	0.89	0.89	0.89	1053
Weighted Avg	0.70	7.00	7.00	1053

Table 2: Classification report summa	rv
--------------------------------------	----

Accuracy	Precision	Recall	F1 Score	AUC
0.9790	0.8968	0.8800	0.8870	0.8920

The results in Table 2 show that the Gradient Boosting Classifier (GBC) achieved an accuracy of 97.9% on the original dataset, with a high F1 score of 88.7%. This indicates that the model performed well in predicting water potability based on the nine original parameters. However, the dataset lacked information on microbial and chemical contaminants, which are critical for a comprehensive assessment of water quality.

The Figure 1 shows that model achieved 399 true negatives and 413 true positives, with 113 false positive and 113 false negative, demonstrating it still makes errors.





To address the limitations of the original dataset, we augmented it with additional parameters, including microbial contaminants (Escherichia coli, Salmonella, Vibrio cholerae, etc.) and chemical contaminants (lead, arsenic, nitrate, etc.). The augmented dataset now included 20 parameters, providing a more comprehensive view of water quality. The performance of the Gradient Boosting Classifier (GBC) on the augmented dataset is summarized Table 3:

Table 3: Classification report				
Parameter	Precision	Recall	F1-Score	Support
0	1.00	1.00	0.98	541
1	1.00	1.00	0.99	512
Accuracy			0.99	1053
Macro Avg	1.00	1.00	1.00	1053
Weighted Avg	1.00	1.00	1.00	1053

Table 4: Classification report summary	
--	--

Accuracy	Precision	Recall	F1 Score	AUC
0.9980	0.9970	0.9970	0.9970	0.9970

The results in Table 4 demonstrate a significant improvement in model performance after augmenting the dataset with additional microbial and chemical parameters. The Gradient Boosting Classifier (GBC) achieved an accuracy of 99.8%, with a nearperfect F1 score of 99.7%. This indicates that the model is highly effective in detecting both microbial and chemical contaminants, making it a robust tool for water quality monitoring.



#### Confusion Matrix

Figure 2: Confusion matrix for the Gradient Boosting Classifier (GBC) on the augmented dataset.

The model achieved 512 true negatives and 541 true positives, with no false positive and no false negative, demonstrating its high accuracy and reliability.

## Experiment II: Feature importance: Original **Dataset vs. Augmented Dataset**

To understand which parameters contributed most to the model's predictions, we conducted a feature importance analysis for both the original and augmented datasets.

# **Feature Importance for Original Dataset**

The feature importance analysis for the original dataset revealed that pH, hardness, and trihalomethanes were

the most significant predictors of water potability. These parameters are commonly associated with water quality but do not account for microbial or chemical contamination.

The results Figure 3 show that pH is the most critical feature influencing the model's predictions, followed by sulfate, hardness, chloramines, and solids. Features like organic carbon, trihalomethanes, conductivity, and turbidity have relatively lower importance.



Figure 3 Graph showing feature importance of parameters in the original dataset



Feature Importance (Gradient Boosting Classifier)

Figure 4: Graph showing feature importance of new parameters in the new dataset

### Feature importance for augmented dataset

The feature importance analysis for the augmented dataset highlighted the significance of the new parameters, particularly Escherichia coli, Salmonella, and lead. These contaminants were found to be critical predictors of water potability, underscoring the importance of including them in the dataset.

Figure 4 shows the new microbial and chemical contaminants sulfate  $(SO_4^{2-})$  (0.0685), copper (Cu) (0.0412), iron (0.0569), lead (Pb) (Fe) (0.0413), chlorine residues (THMs) (0.0423),and aluminum (Al) (0.0452) emerging as very  $(NO_2^-)$ important features, while nitrate

(0.0391), arsenic (As) (0.0372), and Escherichia coli (0.0304) are moderately important. These parameters significantly improve the model's ability to evaluate both chemical and biological risks, ensuring a comprehensive assessment of water quality.

## onclusion

This study demonstrates the effectiveness of expanding the water quality dataset with additional microbial and chemical parameters. combined The augmented dataset. with SMOTE+Tomek balancing, significantly improved the performance of the Gradient Boosting Classifier (GBC), achieving an accuracy of 99.8% and an F1 score of 99.7%. The feature importance analysis confirmed the critical role of the new parameters, particularly Escherichia coli, Salmonella, and lead, in predicting water potability. These findings underscore the importance of comprehensive datasets and advanced machine learning techniques for water quality monitoring in resource-limited settings. Future work should focus on integrating real-time sensor data and validating the model in real-world scenarios to further enhance its applicability and impact.

#### References

- [1] Abolade, O. A., Adewumi, M. O. & Oyedele, O. O. (2024). Assessment of bacteriological quality of sachet water in Ibadan, Nigeria. J. of Envtal Sci. and Techn., 7(2), 124-131. https://doi.org/10.1023/A:1020513008097
- [2] Udoh, A., Akanbi, B. & Essien, N. (2021). Microbial and chemical contamination in sachet water products in Nigeria. Journal of **Environmental** Applied Sciences and Management, 15(3), 41-50. https://doi.org/10.4314/jasem.v15i3.6
- Birhan, M., Tadesse, A. & Assefa, G. (2023). [3] Waterborne diseases in Africa: Trends and public health impact. African Journal of Public Health Research, 12(4), 201–213. https://doi.org/10.11648/j.ajphr.20230402.12
- [4] Li, Y., Zhao, L., Chen, W. & Hu, G. (2024). Heavy metal contamination and public health risks in sachet water. Environmental Science and Pollution Research, 25(3), 2380-2391. https://doi.org/10.1007/s11356-023-27000-x

- Grizzetti, B., Bouraoui, F. & Aloe, A. (2024). [5] Assessment of contaminants in surface and drinking waters across Nigeria. Journal of Environmental Quality, 32(2), 437–446. https://doi.org/10.2134/jeq2023.05.0123
- Li, X., Zhao, J. & Sun, Q. (2023). Enzyme-based [6] biosensors for heavy metal detection in drinking water. Sensors and Actuators B: Chemical, 132471. 376, https://doi.org/10.1016/j.snb.2022.132471
- NCDC (2022). Cholera Outbreak in Nigeria [7] Linked to Contaminated Water Sources. Nigeria Center for Disease Control (NCDC) Annual Report.
- [8] WHO/IWA (2017). Global Status Report on Water Safety Plans: A Review of Proactive Risk Assessment and Risk Management Practices to Ensure the Safety of Drinking-World Health Organization/The Water. International Water Association Press.
- [9] UNICEF (2023). The Impact of Waterborne Diseases on Child Mortality in Nigeria. UNICEF Annual Health Report.
- [10] Fakayode, S. O. & Onianwa, P. C. (2002). Heavy metals contamination of groundwater resources in Ibadan, Nigeria. Environmental Geochemistry and Health, 24(3), 257-266. https://doi.org/10.1023/A:1020513008097
- [11] Akanbi, B. O., Akinola, O. S. & Okonko, I. O. (2020). Microbial contamination of sachet drinking water in Kwara State, Nigeria. J. of Environmental Sci. and Techn., 13(1), 34–40. https://doi.org/10.2134/jeq2023.05.0123
- [12] Olanrewaju, R. (2021). Evaluation of microbial water quality assessment methods in Nigeria. African J. of Microbio. Res., 15(5), 235-245. https://doi.org/10.5897/AJMR2020.10053
- [13] WHO (2017) Guidelines for Drinking-Water Quality: Fourth Edition Incorporating the World First Addendum. Health Organization, Geneva. Available from: https://www.who.int/water\_sanitation\_hea lth/publications/drinking-water-qualityguidelines-4-including-1st-addendum/en/.
- [14] Zhao, L., Liu, Y. & Chen, H. (2024). Predicting water quality changes using machine learning: An approach using ensemble models for dissolved oxygen, pH, and WOI. 113488. Environmental Research, 217, https://doi.org/10.1016/j.envres.2023.113488
- [15] Zhang, T. & Liu, P. (2023). AI-driven forecasting in water quality management: A review. J. of Water Resources Planning and Mgt., 149(3), 04022097. https://doi.org/10.1061/(ASCE)WR.1943-

5452.0001532

[16] Wang, M., Li, X. & Zhou, S. (2024). Real-time water quality monitoring using sensor networks and remote sensing technologies. IEEE Sensors Journal, 24(2), 187-200. https://doi.org/10.1109/JSEN.2023.3320456

https://lafiascijournals.org.ng/index.php/fscproceedings 
FULafia-FSC Conference Proceedings, 2025 161

- [17] Chen, X., Liu, Q., & Zhang, H. (2024). Prediction of water quality index using machine learning models: A case study in China. *Environmental Monitoring and Assessment*, 196(1), 35–49. https://doi.org/10.1007/s10661-023-10994-5
- [18] Li, Y., Zhao, L., Chen, W. & Hu, G. (2024). Heavy metal contamination and public health risks in sachet water. *Environmental Science* and Pollution Research, 25(3), 2380–2391. https://doi.org/10.1007/s11356-023-27000-x
- [19] Inok, Arit & Lawal, Basira Kankia & Akpan, Mary & Labaran, Kamilu & Ndem, Ekpedeme & Ohabunwa, Unoma & Tikare, Olubukola & Idris Ibrahim, Umar & Amorha, Kosisochi & Enevi. Kpokiri. (2021).Microbial contamination of packaged drinking water in Nigeria. Tropical Medicine & International Health. 26. 10.1111/tmi.13672. Stolper, R., Johnson, H. & Karp, A. (2020). Real-time monitoring challenges in water quality assessment: Current methods and innovations. International Journal of Environmental Research and Public Health, 17(5), 1524. https://doi.org/10.3390/ijerph17051524
- [20] Stolper, R., Johnson, H. & Karp, A. (2020). Realtime monitoring challenges in water quality assessment: Current methods and innovations. *International Journal of Environmental Research and Public Health*, 17(5), 1524. https://doi.org/10.3390/ijerph17051524
- [21] Edegbene, Ovie & Yandev, Doowuese & Omotehinwa, Temidayo Oluwatosin & Zakari, Hajara & Andy, Blessing & Ujoh, John. (2024). Exploring the Prevalence of Microorganisms in Selected Water Sources in Benue South of Nigeria using Selected Environmental Factors and Microbial Counts. International Journal of Pathogen Research, 13, 36-45. DOI: 10.9734/ijpr/2024/v13i3284.

- [22] Haghiabi, A. H., Nasrolahi, A. H. & Parsaie, A. (2018). Water quality prediction using machine learning methods. *Water Quality Research Journal*, 53(1), 3–13.
- [23] Ahmed, U., Mumtaz, R., Anwar, H. & Shah, A. A. (2019). Efficient water quality prediction using supervised machine learning. *Water*, 11(11), 2210. https://doi.org/10.3390/w11112210
- [24] El Bilali, A. & Taleb, A. (2020). Prediction of irrigation water quality parameters using machine learning models in a semi-arid environment. *Journal of Saudi Society of Agricultural Sciences*, 19(7), 439–451. https://doi.org/10.1016/j.jssas.2019.05.002
- [25] Aldhyani, T. H. H., Al-Yaari, M., Alkahtani, H. & Maashi, M. (2020). Water quality prediction using artificial intelligence algorithms. *Applied Bionics and Biomechanics*, 2020, 6659314. https://doi.org/10.1155/2020/6659314
- [26] Lu, H. & Ma, X. (2020). Hybrid decision treebased machine learning models for short-term water quality prediction. *Environmental Modelling & Software*, 124, 104604. https://doi.org/10.1016/j.envsoft.2019.104604
- [27] Dritsas, E. & Trigka, M. (2023). Water quality classification using machine learning models. *Water*, 15(2), 345. https://doi.org/10.3390/w15020345
- [28] Wang, H., Zhang, L. & Wang, J. (2023). Water quality prediction using machine learning models: A case study of Nainital Lake. *Environmental Science and Pollution Research*, 30(1), 1–12. https://doi.org/10.1007/s11356-022-21124-8