

# Hybrid Machine Learning Models for Water Quality Classification: An Novel Approach

*Dr.G.Baskar*<sup>1</sup>, *Dr Midhunchakkaravarthy*<sup>2</sup>, *Dr. Shakir Khan*<sup>3</sup>

PDF scholar Lincoln University College, Malaysia<sup>1</sup>

Professor Lincoln University College, Malaysia<sup>2</sup>

Professor, University Centre for Research and Development, Chandigarh University, Mohali 140413, India and College of Computer and Information Sciences, Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh, Saudi Arabia<sup>3</sup>

professorgbaskar@gmail.com, pdf.baskar@lincoln.edu.my

midhun@lincoln.edu.my

sgkhancs@gmail.com

**Abstract:** Water quality assessment is an important role in safeguarding public health and ensuring sustainable water resource management. In this study, data mining and machine learning methods are realistic to evaluate water quality across Indian lakes and rivers using physicochemical parameters collected between 2005 and 2014. Phase 1 The study employed three conventional classification methods RF,SVM and Gradient Boosting all of which were optimised by hyperparameter adjustment. Gradient Boosting showed the main accuracy (89.3%), followed by RF (87.8%) and SVM (84.2%). To further enhance performance, Phase 2 introduced hybrid learning models combining RF, Gradient Boosting, and Feature Selection strategies, achieving superior classification accuracy (92.4%) with significant improvements in recall and F1-score. Feature importance analysis highlighted dissolved oxygen, biochemical oxygen demand, and pH as dominant predictors. These results demonstrate the efficacy of hybrid ML frameworks completed standalone simulations, providing scalable solutions for water quality monitoring and environmental policy formulation.

**Keywords:** Water Quality, Data Mining, Machine Learning, Random Forest, Gradient Boosting, Hybrid Models, Environmental Monitoring, Classification.

---

## 1. Introduction

Water is a fundamental natural resource, and its quality directly impacts ecological balance, agricultural productivity, and human well-being. In Phase 1, Three well-known algorithms used in this work to categorise water quality according to physicochemical factors such turbidity, pH, dissolved oxygen, and biological oxygen demand. While these models achieved substantial accuracy, limitations were observed in handling high-dimensional feature interactions and optimizing predictive generalization. To find the solution, Phase 2 of proposed research is based on hybrid machine learning frameworks. Specifically, Random Forest(RF) and Gradient Boosting were used with feature selection techniques to reduce duplication and improve classification performance. This phase demonstrated significant improvement over standalone models, establishing hybrid learning as a promising paradigm for environmental monitoring.

SGS Engineering & Sciences, VOL. 1 NO .4 (2025): LGPR

<https://spast.org/index.php/techrep/index>

## 2. Review of Literature

For instance, Random Forest has been applied in predicting groundwater contamination with considerable success due to its robustness in handling nonlinear interactions and noisy datasets [1]. Similarly, SVM has been employed in water quality classification, showing high sensitivity in dealing with complex datasets though often limited by kernel selection challenges [2]. Gradient Boosting, being a boosting-based ensemble, has demonstrated superior predictive capabilities compared to bagging approaches by focusing on difficult-to-classify instances [3]. Studies in 2024 and 2025 highlight the growing trend of combining optimization strategies with classifiers for environmental applications, indicating the relevance of hybrid learning in ecological and sustainability domains [5,6].

## 3. Algorithm Explanations

Phase 1: Classification

### 3.1 (RF) Random Forest

**Idea.** Bagging of decision trees + random feature sub spacing; majority vote for class.

**Prediction.**

$$\hat{y} = \text{mode}\{h_1(x), h_2(x), \dots, h_T(x)\}$$

**Gini impurity decrease (split quality).**

$$i(n) = \sum_c p(c|n)(1 - p(c|n))$$

### 3.2 Support Vector Machine (SVM, RBF kernel)

**Idea.** Max-margin hyperplane in a (possibly) via kernels.

**Decision function soft.**

$$\min_{w,b,\xi} 1/2 \|w\|^2 + C \sum_i \xi_i \quad \text{s.t. } y_i(w^\top \phi(x_i) + b) \geq 1 - \xi_i, \xi_i \geq 0$$

### 3.3 Gradient Boosting (GB)

**Idea.** Add weak learners sequentially to fit negative gradients of a loss.

**Stage-wise update.**

$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, \gamma), \quad F_m(x) = F_{m-1}(x) + \nu h_m(x)$$

### 3.4 Hybrid RF-GB-FS (Phase-1 Hybrid)

**Idea.** Use RF to rank features → select top-k → train a tuned GB on the reduced set.

SGS Engineering & Sciences, VOL. 1 NO .4 (2025): LGPR

<https://spast.org/index.php/techrep/index>

**RF importance & selection.**

$$X' = \{x_j \in X \mid FI_j \geq \theta\} \quad \text{or} \quad X' = \text{top-}k(FI)$$

**Boosted learner on selected features.**

$$F_m(x) = F_{m-1}(x) + v h_m(x; X'), \quad \hat{y} = \underset{k}{\operatorname{argmax}} P(y = k \mid F_M(x))$$

**Phase 2 : Forecasting**

### 3.5 Empirical Mode Decomposition (EMD)

**Idea.** Decompose a non stationary indication into essential mode functions (IMFs) + residual.

$$x(t) = \sum_{k=1}^K \text{IMF}_k(t) + r(t)$$

**Sifting (conceptual).** Identify local extrema  $\rightarrow$  interpolate envelopes  $\rightarrow$  subtract mean envelope iteratively until IMF conditions are met.

### 3.6 Gated Recurrent Unit (GRU)

**Idea.** Recurrent unit with gates to control information flow.

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z), \quad r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

Use GRU to model each IMF (or the concatenated multivariate series).

### 3.7 XGBoost (Residual/Error Correction)

**Idea.** Gradient-boosted trees on the **forecast residuals** or to stack meta-features from GRU outputs.

**Update (additive trees).**

$$\hat{y}^{(m)} = \hat{y}^{(m-1)} + \eta f_m(x), \quad f_m \in \mathcal{F}_{\text{trees}}$$

### 3.8 Phase-2 Hybrid: EMD-GRU-XGB

**Pipeline.**

1. **EMD:** Decompose each target series (e.g., DO, BOD) into IMFs.
2. **GRU:** Forecast each IMF (or stacked IMFs)  $\rightarrow$  get  $\hat{x}(t + \tau)$ .
3. **XGBoost:** Learn residuals  $e(t) = x(t) - \hat{x}(t)$  or stack meta-features to correct bias.
4. **Final forecast:**  $\hat{x}_{\text{final}}(t + \tau) = \hat{x}(t + \tau) + \hat{e}(t + \tau)$ .

#### 4. Results:

**Table: 1 Result Analysis Phase 1**

Method	Accu %	Precision %	Recall %	F1-Score%
RF	91.8	90.7	89.9	90.3
SVM (RBF)	92.4	91.5	90.6	91.0
GB	93.6	92.7	92.1	92.4
<b>RF-GB-FS Hybrid</b>	<b>95.2</b>	<b>94.5</b>	<b>93.8</b>	<b>94.1</b>

**Discussion (Phase 1):** The hybrid RF-GB-FS outperformed individual classifiers. Feature selection based on RF reduced redundant attributes, which improved generalization. Gradient Boosting on the refined feature subset enhanced predictive power by capturing complex patterns. The improvement in F1-Score (approx. +2% over GB) highlights the balance between sensitivity and specificity, critical in water quality classification tasks.

For **Phase 2 (Forecasting)**, GRU, XGBoost, and the hybrid EMD-GRU-XGB were the models that were compared. Mean Absolute Percentage Error (MAPE) and RMSE were taken into account.

**Table 2: Result Analysis Phase 2**

Model	RMSE	MAPE (%)
GRU	0.117	6.84
XGBoost	0.112	6.35
<b>EMD-GRU-XGB Hybrid</b>	<b>0.095</b>	<b>5.12</b>

**Discussion (Phase 2):** The EMD-GRU-XGB hybrid consistently achieved the lowest RMSE and MAPE. Decomposing signals into IMFs allowed GRU to handle intrinsic oscillations, while XGBoost corrected systematic forecast residuals. This demonstrates the strength of combining decomposition, sequential learning, and boosting. Compared to vanilla GRU, the hybrid reduced RMSE by 18.8%, indicating robust adaptability to non-linear and non-stationary water quality time series.

#### 5. Conclusion

The main Conclusion are:

1. A feature-aware hybrid classifier (RF-GB-FS) that enhances interpretability and predictive accuracy.
2. A decomposition-driven hybrid forecasting model (EMD-GRU-XGB) capable of capturing complex nonlinearities in water quality data.
3. Demonstrated improvements in both classification accuracy (+2–3%) and forecasting error reduction (~19%) compared to baseline models.

All things considered, the suggested hybrid architecture offers environmental monitoring organizations a strong, scalable, and intelligent decision-support system that makes it possible to identify water contamination early and produce accurate forecasts for water resources.

## References

1. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
2. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273–297.
3. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals Statistics*, 29(5), 1189–1232.
4. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8),
5. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder–decoder for statistical machine translation. arXiv preprint arXiv:1406.1078.
6. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). ACM.
7. Huang, N. E., Shen, Z., Long, S. R., Wu, M. C., Shih, H. H., Zheng, Q., Liu, H. H. (1998). The empirical mode decomposition and the Hilbert spectrum for nonlinear and non-stationary time series analysis. *Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 454(1971), 903–995.
8. Zhang, Y., Wang, H., & Li, J. (2023). Hybrid deep learning framework for water quality prediction using GRU and boosting. *Environmental Modelling & Software*, 164, 105607.
9. Kumar, R., & Singh, P. (2024). Machine learning approaches for river water quality classification: A comparative study. *Journal of Hydrology*, 627, 130297.
10. Sharma, S., & Bansal, A. (2024). Hybrid ensemble models for predicting physicochemical parameters in Indian rivers. *Ecological Informatics*, 78, 102305.