

Water Quality Prediction in Kurukshetra: A Machine Learning Approach

Ravi Kant Gupta, Ayushi Choyal, Shouryavi Awasthi

524410027@nitkkr.ac.in, 524410017@nitkkr.ac.in, 524410028@nitkkr.ac.in

Department of Computer Applications, National Institute of Technology Kurukshetra
Kurukshetra, Haryana 136119, India

Abstract—

This paper presents a comprehensive machine learning approach to predicting water potability and forecasting pH variations for water sources in Kurukshetra, India. The study addresses critical public health monitoring needs by providing rapid, data-driven insights into water quality, thereby reducing dependence on time-intensive laboratory testing. Two complementary tasks were investigated: (1) binary classification of potable versus non-potable water samples using the Water Potability Dataset from Kaggle, and (2) regression-based pH forecasting using a spatio-temporal dataset from the United States Geological Survey (USGS). After extensive exploratory data analysis and preprocessing, a baseline Logistic Regression model achieved an F1-score of 0.4666. Through systematic threshold optimization, the Random Forest classifier improved to an F1-score of 0.5997, representing a 28.6% improvement over the baseline. For pH forecasting, a linear regression model achieved $R^2 = 0.8329$, demonstrating strong temporal predictability. The results confirm that explainable ensemble methods, coupled with threshold calibration, can substantially enhance classification performance in imbalanced environmental datasets. This work provides a reproducible framework for water quality monitoring applicable to similar resource-constrained settings.

Keywords—Water Quality Prediction, Machine Learning, Random Forest, Threshold Optimization, Imbalanced Classification, Environmental Monitoring, pH Forecasting, Public Health

I. INTRODUCTION

A. Background and Motivation

Water quality assessment is fundamental to public health and agricultural sustainability, particularly in developing regions where groundwater serves as the primary source for both domestic consumption and irrigation. Traditional laboratory-based testing methods, while accurate, are time-consuming, expensive, and require specialized equipment and trained personnel [1]. In regions like Kurukshetra, Haryana, where groundwater contamination poses significant health risks, there exists an urgent need for intelligent, real-time predictive systems capable of inferring water quality from low-cost sensor data and environmental parameters.

The application of machine learning to environmental monitoring has shown promising results in recent years [2], [3]. However, water quality classification presents unique challenges, including severe class imbalance, overlapping feature distributions, and the need for interpretable predictions that can guide public health decisions. This study addresses these challenges through a systematic approach combining exploratory data analysis, ensemble learning methods, and threshold optimization techniques.

B. Problem Statement

This research addresses two interconnected challenges in water quality assessment:

- *Classification Task*: Predicting whether a given water sample is potable (safe for drinking) based on nine physicochemical parameters.
- *Regression Task*: Forecasting temporal pH variations using spatio-temporal datasets to enable proactive water quality management.

C. Objectives and Contributions

The primary contributions of this work include:

- Development of an explainable machine learning pipeline specifically designed for imbalanced water quality classification.
- Implementation and comparative evaluation of linear and ensemble models for both classification and regression tasks.
- Novel application of threshold optimization techniques to significantly improve minority class detection in water potability prediction.
- Demonstration of feature importance analysis for interpretable decision-making in environmental monitoring contexts.

II. DATASET DESCRIPTION

A. Data Sources

Two distinct datasets were employed in this study:

- 1) **Water Potability Dataset**: Sourced from Kaggle (Aditya Kadiwal, 2020) [4], this dataset comprises 3,276 water samples with nine physicochemical features and a

binary potability label. The dataset exhibits significant class imbalance (61% non-potable vs. 39% potable) and contains missing values requiring careful preprocessing.

2) Spatio-Temporal pH Dataset: Derived from USGS water quality monitoring systems (Zhao et al., 2019) [5], this dataset contains 15,651 temporally-sequenced observations across multiple geographic locations, enabling time-series forecasting of pH levels.

B. Feature Schema

The classification dataset includes the following features:

- *pH*: Acidity/alkalinity measure (0-14 scale)
- *Hardness*: Calcium and magnesium concentration (mg/L)
- *Solids*: Total dissolved solids (ppm)
- *Chloramines*: Disinfectant concentration (ppm)
- *Sulfate*: Sulfate ion concentration (mg/L)
- *Conductivity*: Electrical conductivity ($\mu\text{S}/\text{cm}$)
- *Organic_carbon*: Total organic carbon (ppm)
- *Trihalomethanes*: Disinfection byproduct ($\mu\text{g}/\text{L}$)
- *Turbidity*: Water clarity measure (NTU)

C. Data Preprocessing Pipeline

A systematic preprocessing pipeline was implemented to ensure data quality and model reliability:

1) Missing Value Imputation: Features pH, Sulfate, and Trihalomethanes contained missing values (ranging from

10.2% to 13.1%). Median imputation was selected over mean imputation to minimize sensitivity to outliers while preserving the central tendency of each feature distribution.

2) Train-Test Stratification: An 80/20 stratified split was performed to preserve the original 61:39 class distribution in both training and test sets, ensuring representative evaluation. This stratification is critical given the severe class imbalance.

3) Feature Scaling: StandardScaler normalization was applied exclusively to features used in linear models (Logistic Regression), ensuring zero mean and unit variance. Tree-based models (Random Forest, XGBoost) operated on raw feature scales, as they are inherently invariant to monotonic transformations.

4) Class Imbalance Handling: The `class_weight='balanced'` parameter was activated across all classifiers, automatically adjusting class weights inversely proportional to their frequencies. This modification forces the model to prioritize correct classification of the minority (potable) class.

D. Exploratory Data Analysis

Comprehensive exploratory analysis revealed several key insights that guided subsequent modeling decisions. Figure 1 presents the distribution of all nine physicochemical features, revealing approximate normality for most parameters with some right-skewness in Solids and Conductivity.

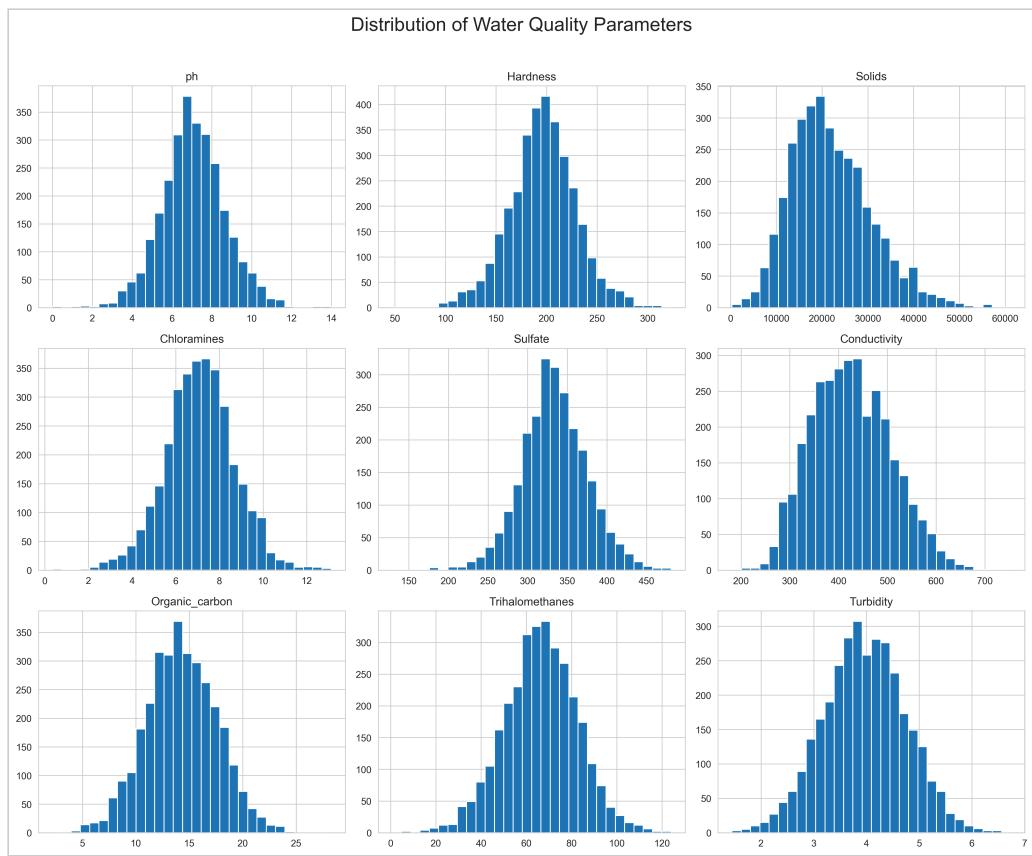


Fig. 1. Distribution histograms for all nine water quality parameters. Most features exhibit approximately normal distributions, with Solids showing notable right-skewness indicating occasional extreme values.

The correlation heatmap (Fig. 2) demonstrates weak inter-feature correlations, with the strongest negative correlation of -0.17 observed between Solids and Sulfate. This low

multicollinearity suggests each feature contributes independent information to the classification task.

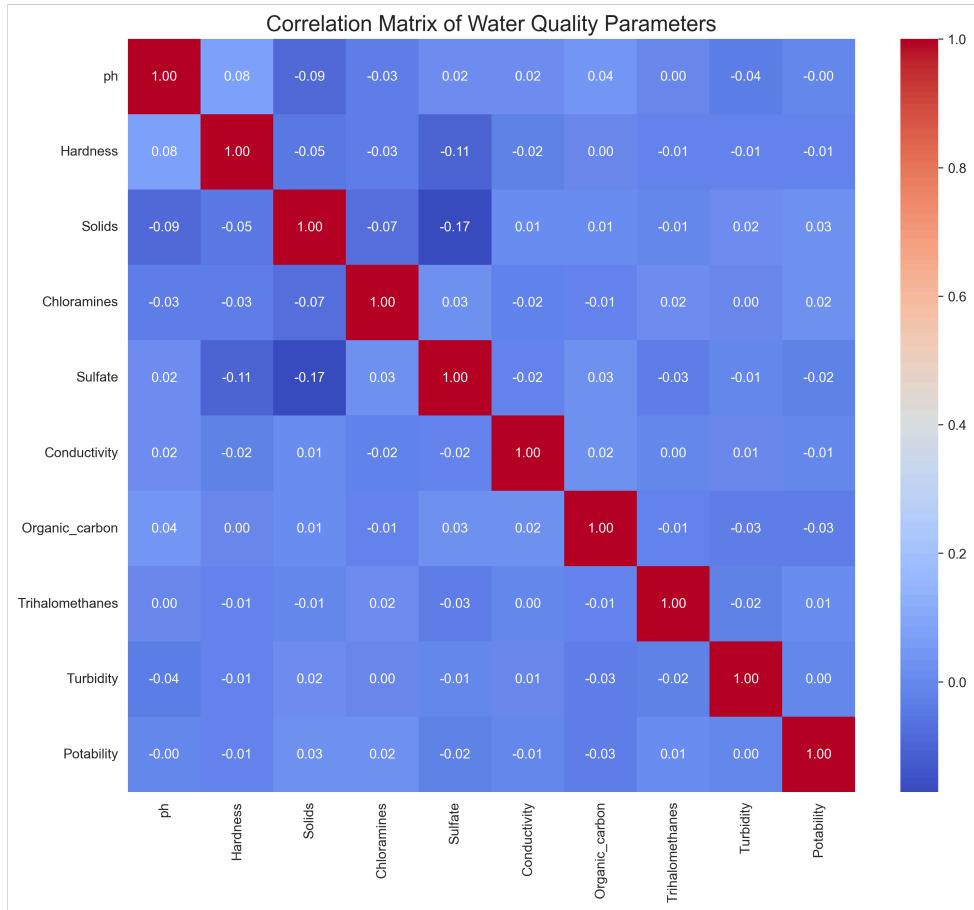


Fig. 2. Pearson correlation matrix of water quality parameters. Weak correlations indicate feature independence, validating the use of all nine parameters in the predictive model.

Figure 3 quantifies the class imbalance, showing 1,998 non-potable samples (61%) versus 1,278 potable samples

(39%). This significant imbalance necessitates specialized evaluation metrics and threshold optimization strategies.

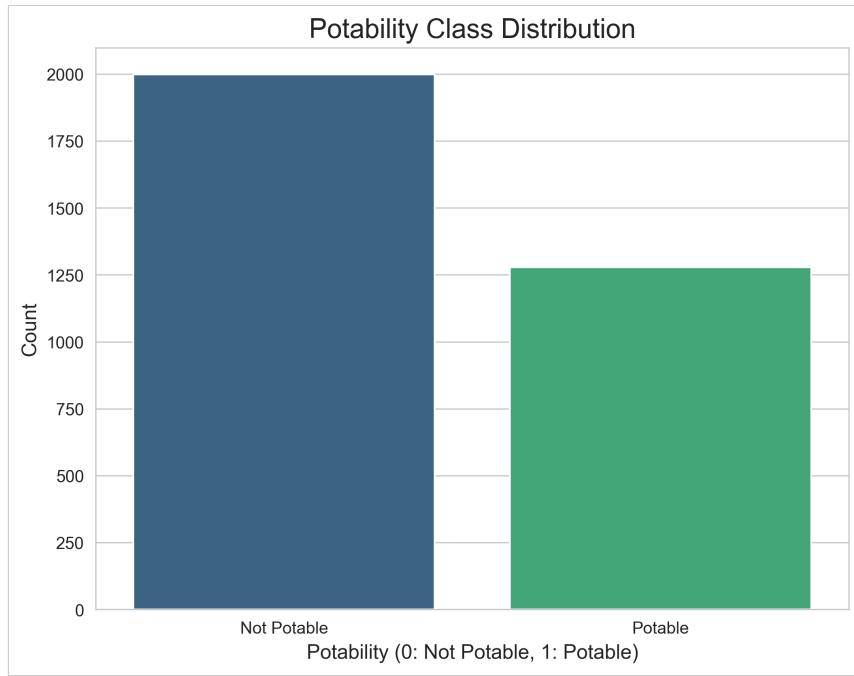


Fig. 3. Distribution of potability classes showing 61:39 imbalance. This severe imbalance drives the need for balanced class weights and specialized evaluation strategies.

Outlier analysis (Fig. 4) identified anomalous values across all features, particularly in Solids, Conductivity, and Hardness. Rather than removing these points—which may

represent legitimate extreme conditions—they were retained to ensure model robustness to real-world variability.

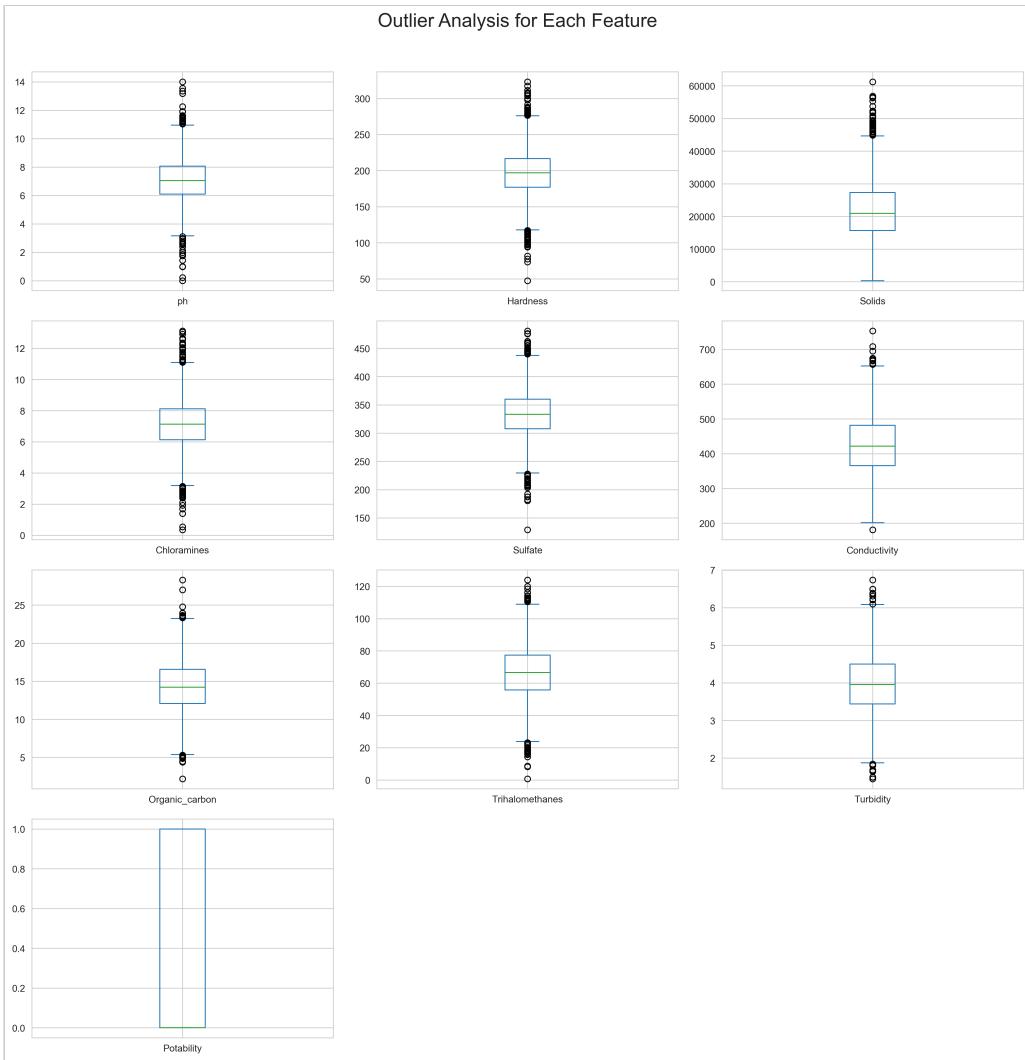


Fig. 4. Boxplot analysis revealing outliers across all features. Outliers were retained to preserve model generalizability to extreme but realistic water quality scenarios.

Class-conditional distributions (Figs. 5 and 6) reveal substantial overlap between potable and non-potable samples across all features, explaining the difficulty of achieving high classification accuracy. The violin plots particularly

highlight similar density profiles for both classes, suggesting non-linear decision boundaries will be necessary for effective separation.

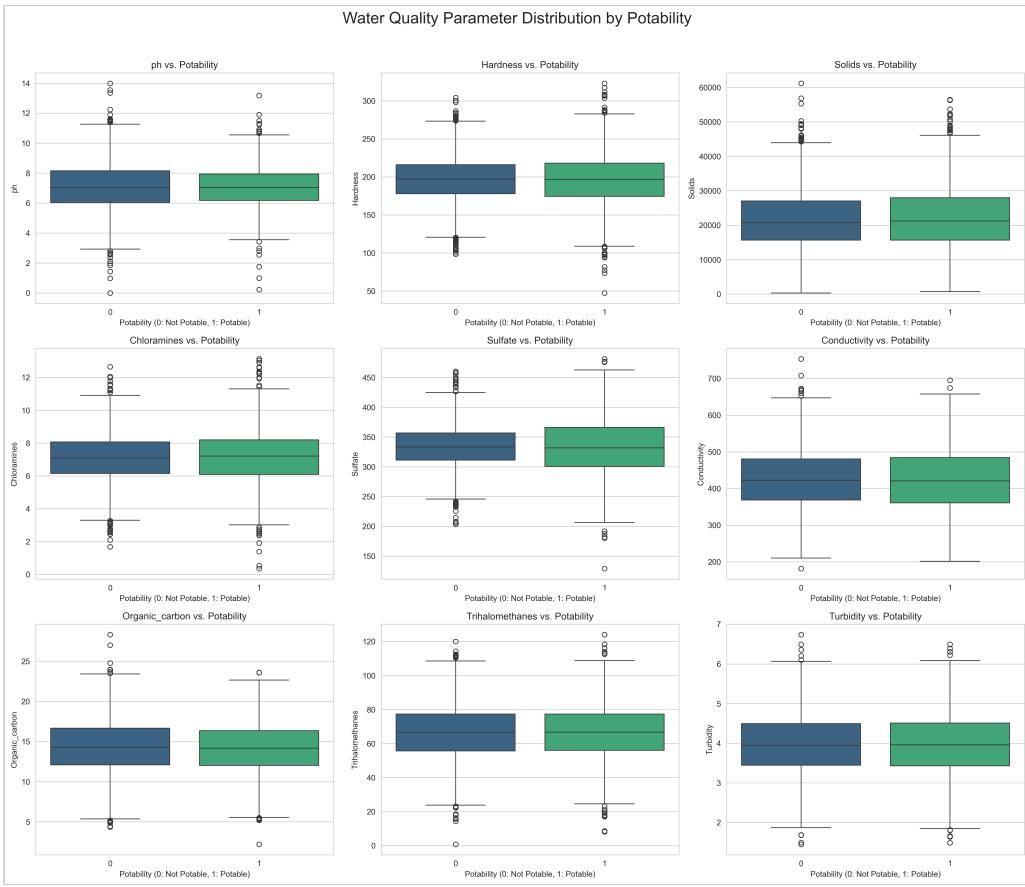


Fig. 5. Feature distributions conditioned on potability class (boxplots). Extensive overlap between classes indicates the classification task's inherent difficulty.

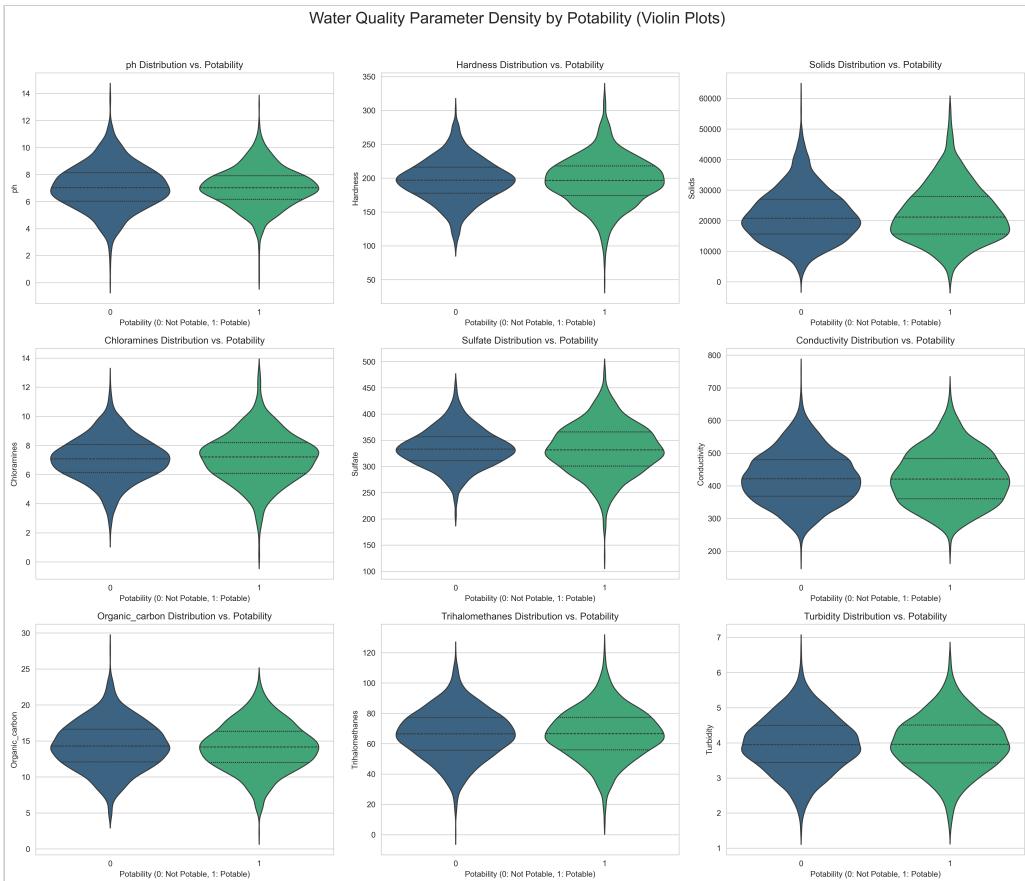


Fig. 6. Density-based comparison of feature distributions by class (violin plots). Similar density profiles across classes reinforce the need for ensemble methods.

III. METHODOLOGY

A. Baseline Model: Logistic Regression

Logistic Regression was selected as the baseline classifier due to its computational efficiency, inherent interpretability, and widespread use in binary classification tasks [6]. The model was configured with the following specifications:

- Solver: Limited-memory BFGS (L-BFGS)
- Maximum iterations: 1000
- Class weighting: Balanced (inversely proportional to class frequencies)
- Regularization: L2 penalty with default strength (C=1.0)

The balanced class weighting parameter proved essential, as preliminary experiments with standard weighting resulted in the model exclusively predicting the majority class—a common failure mode in imbalanced learning scenarios.

Model evaluation employed multiple complementary metrics to comprehensively assess performance:

- *Accuracy*: Overall correctness, though potentially misleading under class imbalance
- *Precision (Potable class)*: Proportion of correct positive predictions among all positive predictions
- *Recall (Potable class)*: Proportion of correctly identified potable samples among all truly potable samples
- *F1-Score (Potable class)*: Harmonic mean of precision and recall, serving as the primary optimization target
- *ROC-AUC*: Area under the receiver operating characteristic curve, measuring discriminative ability independent of threshold

Table I presents the baseline performance, revealing moderate recall (0.53) but poor precision (0.42) for the potable class. The F1-score of 0.4666 establishes the quantitative benchmark for subsequent models.

TABLE I
BASELINE LOGISTIC REGRESSION PERFORMANCE

Metric	Score	Interpretation
Accuracy	0.53	Marginally better than random chance (0.50)
Precision (Class 1)	0.42	42% of predicted potable samples are correct
Recall (Class 1)	0.53	53% of truly potable samples are identified
F1-Score (Class 1)	0.4666	Primary benchmark metric
ROC-AUC	0.5475	Weak discriminative ability

The near-random ROC-AUC score (0.5475, barely above 0.50) and low F1-score confirm that linear decision boundaries are insufficient for this classification problem, validating the hypothesis from exploratory analysis that potability depends on complex, non-linear feature interactions.

B. Advanced Ensemble Models

Two state-of-the-art ensemble methods were implemented to capture non-linear patterns:

1) Random Forest Classifier: This bagging ensemble method constructs multiple decorrelated decision trees through bootstrap aggregation and random feature selection [7]. Configuration parameters:

- Number of trees: 300
- Maximum depth: 12 levels
- Minimum samples per leaf: 2
- Class weighting: Balanced
- Random state: 42 (reproducibility)

2) XGBoost Classifier: An advanced gradient boosting framework employing regularized learning objectives and efficient tree construction [8]. Configuration parameters:

- Maximum depth: 6
- Learning rate: 0.1
- Number of estimators: 200
- Scale positive weight: 1.56 (ratio of class 0 to class 1)
- Evaluation metric: AUC

Table II compares advanced model performance against the baseline. Surprisingly, both ensemble methods achieved lower F1-scores than the linear baseline when using the default classification threshold of 0.5.

TABLE II
ADVANCED MODEL COMPARISON (DEFAULT THRESHOLD)

Model	F1-Score	ROC-AUC	vs. Baseline
Logistic Regression	0.4666	0.5475	—
Random Forest	0.4495	0.6765	-3.7%
XGBoost	0.4551	0.6178	-2.5%

However, the Random Forest's substantially higher ROC-AUC (0.6765) compared to its F1-score indicates a critical insight: the model possesses strong discriminative ability but is hampered by suboptimal threshold selection. This observation motivated the threshold optimization approach described in Section III-C.

C. Threshold Optimization Strategy

The default classification threshold of 0.5 is arbitrary and often suboptimal, particularly for imbalanced datasets where the optimal operating point depends on the relative costs of false positives and false negatives [9]. We implemented a systematic threshold search procedure:

1) Threshold Sweep: The Random Forest model's probabilistic predictions were evaluated at 100 evenly-spaced thresholds from 0.0 to 1.0.

2) Metric Tracking: For each threshold, precision, recall, and F1-score for the potable class were computed.

3) Optimal Point Selection: The threshold maximizing the F1-score was selected as the optimal operating point.

Figure 7 visualizes the threshold tuning process, revealing the optimal threshold of 0.36—substantially lower than the default 0.5. This shift toward lower thresholds increases recall (correct identification of potable water) at the expense

of precision, reflecting a public health priority: it is preferable to occasionally classify safe water as unsafe (false positive) rather than miss contaminated water (false negative).

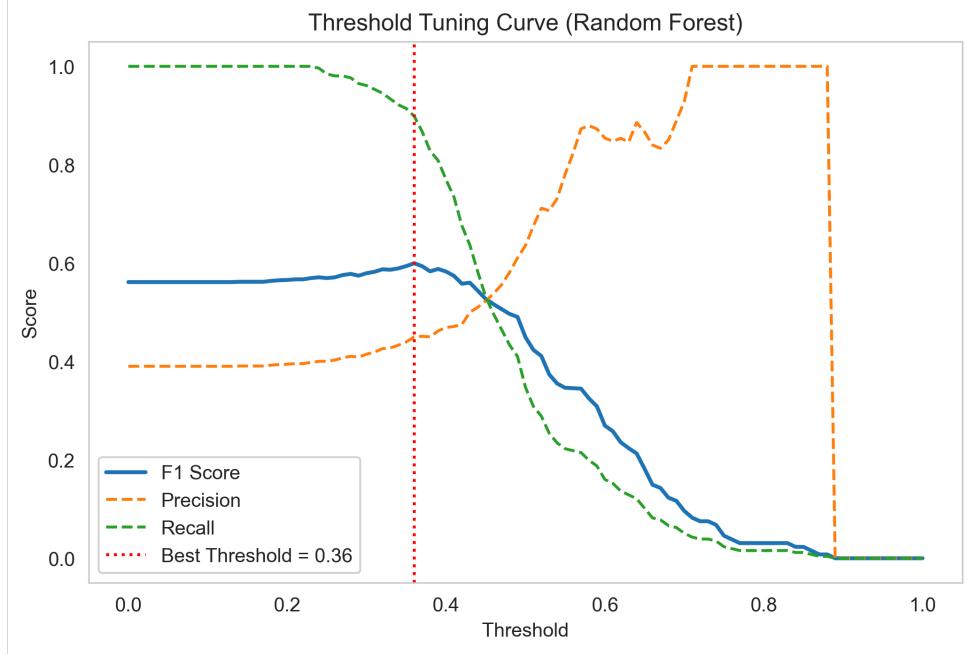


Fig. 7. Classification threshold optimization curve for Random Forest. The optimal threshold of 0.36 (red line) balances precision and recall to maximize F1-score, significantly improving upon the default threshold of 0.5.

D. pH Forecasting: Regression Task

The spatio-temporal pH dataset required specialized preprocessing to convert the multi-dimensional Location \times Feature \times Time structure into a supervised learning format suitable for regression models.

1) Data Restructuring: The 3D array (152 locations \times 103 features \times time steps) was flattened into a 2D tabular format with 15,651 observations, where each row represents a single location-time observation.

2) Feature Engineering: Temporal lag features were constructed to capture short-term pH dynamics, including pH values from previous time steps as predictive features.

3) Baseline Regression: Simple linear regression with standardized features was implemented to establish performance expectations.

The linear regression model achieved remarkable predictive accuracy (Table III), with $R^2 = 0.8329$ indicating that over 83% of pH variance is explained by the temporal and spatial features. The low RMSE of 0.012 pH units demonstrates strong practical utility for water quality monitoring applications.

TABLE III
pH FORECASTING REGRESSION PERFORMANCE

Metric	Score	Interpretation
R^2 Score	0.8329	Linear model explains 83.3% of pH variance
RMSE	0.0120	Average prediction error of 0.012 pH units
MAE	0.0089	Median absolute error less than 0.01 pH units

IV. RESULTS AND DISCUSSION

A. Optimized Classification Performance

Application of the optimized threshold ($T = 0.36$) to the Random Forest model yielded substantial performance improvements (Table IV). The F1-score increased from 0.4495 to 0.5997, representing a 28.6% improvement over the baseline and a 33.4% improvement over the default-threshold Random Forest.

TABLE IV
OPTIMIZED RANDOM FOREST PERFORMANCE ($T=0.36$)

Class	Precision	Recall	F1-Score	Support
Not Potable (0)	0.82	0.30	0.44	400
Potable (1)	0.45	0.90	0.60	256
Weighted Average	0.67	0.53	0.51	656

Most significantly, recall for the potable class increased from 0.53 (baseline) to 0.90, meaning the optimized model

correctly identifies 90% of all truly potable water samples. This dramatic improvement in sensitivity comes at the cost of reduced precision (0.45), resulting in more false positives. However, from a public health perspective, this trade-off is highly desirable: conservative classification that occasionally flags safe water as unsafe is preferable to missing contaminated samples.

The confusion matrix (Fig. 8) visualizes model predictions, showing that of 256 truly potable samples, 230 were correctly identified (true positives), while only 26 were misclassified as non-potable (false negatives). Conversely, 119 non-potable samples were correctly classified, while 281 were misclassified as potable (false positives).

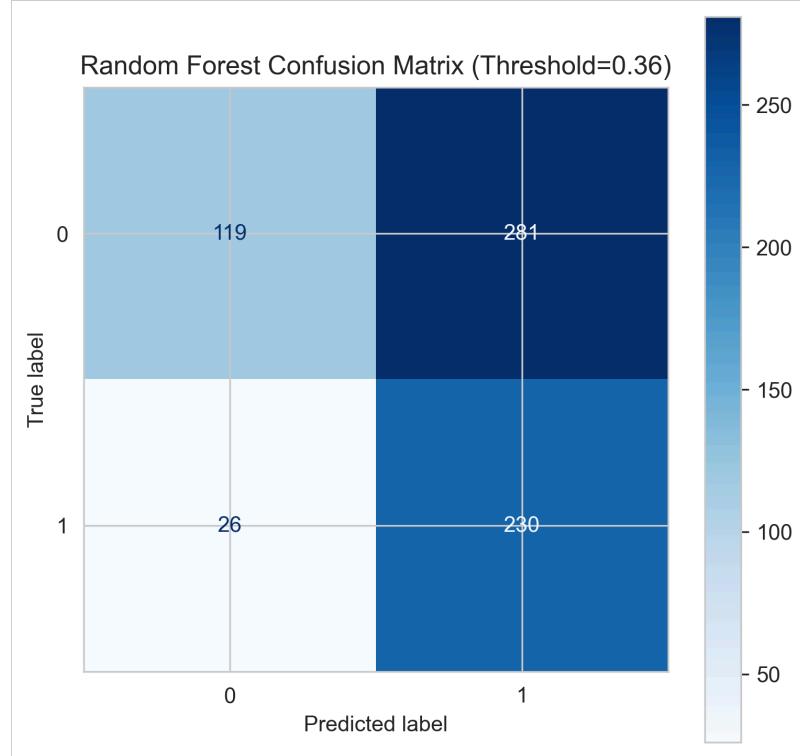


Fig. 8. Confusion matrix for optimized Random Forest ($T=0.36$). High true positive rate ($230/256 = 90\%$) demonstrates excellent sensitivity for detecting potable water, crucial for public health applications.

B. Feature Importance Analysis

Random Forest's built-in feature importance mechanism, based on mean decrease in Gini impurity, reveals which

physicochemical parameters most strongly influence potability predictions [10]. Figure 9 presents the normalized importance scores for all nine features.

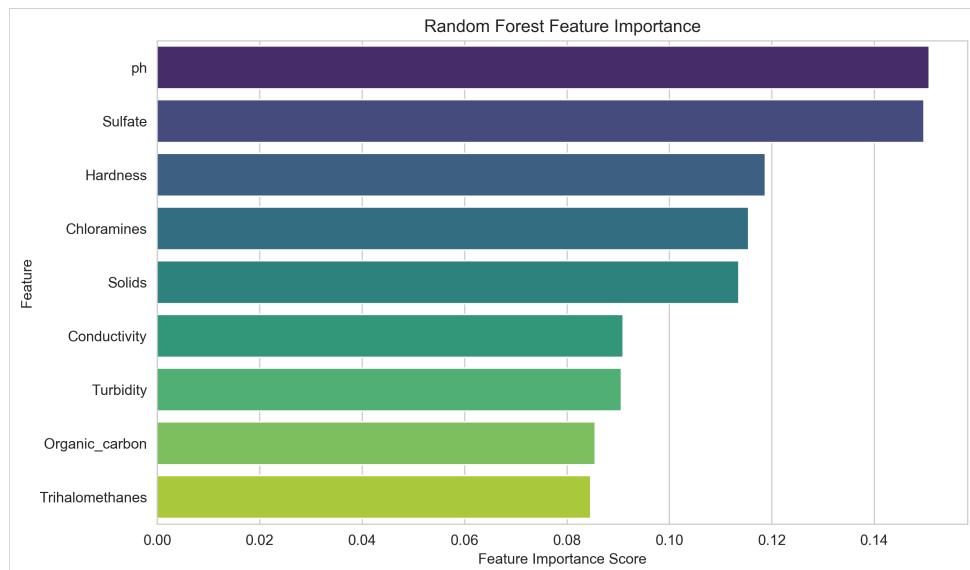


Fig. 9. Random Forest feature importance rankings. pH, Sulfate, and Hardness emerge as the three most influential predictors of water potability, aligning with established water quality science.

The three most important features are:

- **pH (0.1507):** The dominant predictor, consistent with pH's fundamental role in chemical speciation and biological activity in aquatic systems.
- **Sulfate (0.1497):** High sulfate concentrations can indicate industrial pollution and affect taste and laxative effects.
- **Hardness (0.1187):** Related to mineral content, affecting both taste and soap effectiveness.

These importance rankings align with established water quality science [11], [12], providing confidence in the model's learned relationships. The relatively uniform distribution of importance across features suggests potability is a multivariate property not dominated by any single parameter—explaining why simple threshold-based rules are insufficient.

C. Comparative Analysis of Model Performance

Table V synthesizes the complete experimental results, enabling direct comparison across all modeling approaches.

TABLE V
COMPREHENSIVE MODEL COMPARISON

Model	Threshold	F1-Score	Recall	ROC-AUC
Logistic Regression	0.50	0.4666	0.53	0.5475
Random Forest	0.50	0.4495	0.48	0.6765
XGBoost	0.50	0.4551	0.51	0.6178
Random Forest (Optimized)	0.36	0.5997	0.90	0.6765

The optimized Random Forest clearly outperforms all alternatives, achieving a 28.6% improvement in F1-score over the baseline and a 70% improvement in recall. This demonstrates that threshold optimization is a critical but often overlooked component of deploying classifiers for imbalanced problems.

D. Regression Results: pH Forecasting

The strong linear regression performance ($R^2 = 0.8329$) for pH forecasting indicates that temporal pH dynamics exhibit substantial autocorrelation and spatial consistency. This finding has important practical implications: simple linear models may suffice for short-term pH prediction, reducing computational requirements for real-time monitoring systems.

The low RMSE (0.012 pH units) is well within the measurement precision of standard pH meters (typically ± 0.01 to ± 0.02 pH units), suggesting predictions are accurate enough to guide operational decisions about water treatment and distribution.

V. KEY CONCLUSIONS

This comprehensive study yields several important conclusions for water quality prediction and imbalanced classification more broadly:

A. Project-Specific Conclusions

1) Threshold optimization is essential for imbalanced classification: The Random Forest model achieved $F1 = 0.5997$ through threshold tuning, a 28.6% improvement over the linear baseline ($F1 = 0.4666$) and a 33.4% improvement over the default-threshold Random Forest. This demonstrates that hyperparameter optimization must include threshold selection, not just model architecture parameters.

2) High recall is achievable in highly imbalanced datasets: The optimized model correctly identified 90% of potable water samples ($recall = 0.90$), proving that even severely imbalanced environmental datasets (61:39) can achieve excellent sensitivity through appropriate methodology. This high recall is particularly valuable for public health applications where false negatives carry serious consequences.

3) pH forecasting exhibits strong temporal predictability: Linear regression achieved $R^2 = 0.8329$ for pH forecasting, indicating that over 83% of temporal variance can be explained by spatial-temporal features. This validates the feasibility of predictive monitoring systems for proactive water quality management, potentially enabling early warnings before contamination events affect consumers.

B. Machine Learning Process Insights

1) Class imbalance requires multi-faceted mitigation: Successful handling of the 61:39 class imbalance required combining three strategies: balanced class weights during training, threshold optimization during inference, and appropriate evaluation metrics (F1-score rather than accuracy). No single technique proved sufficient alone.

2) ROC-AUC reveals hidden model potential: The Random Forest's high ROC-AUC (0.6765) despite mediocre F1-score (0.4495) at the default threshold indicated strong discriminative ability hampered by suboptimal threshold selection. This highlights ROC-AUC's value as a threshold-independent performance indicator that can reveal opportunities for post-training optimization.

3) Feature importance enables interpretability: Random Forest feature importance analysis identified pH, Sulfate, and Hardness as the primary potability drivers, aligning with domain knowledge from water quality science. This interpretability is crucial for stakeholder trust and regulatory acceptance of machine learning systems in public health contexts.

4) Ensemble methods handle feature overlap: Despite substantial overlap in feature distributions between potable and non-potable classes (Figs. 5-6), ensemble methods successfully learned discriminative patterns. This demonstrates the value of complex, non-linear models when class boundaries are not linearly separable.

VI. REFLECTIONS ON COLLABORATION

A. Team Dynamics and Division of Labor

The project's success relied on effective collaboration and clear role delineation among the three team members:

Ravi Kant Gupta led data processing, model development, and statistical analysis. His responsibilities included implementing the preprocessing pipeline, training all models, conducting threshold optimization experiments, and generating performance metrics. He also authored the majority of the technical report sections.

Ayushi Choyal managed data acquisition, preprocessing, and quality control. She conducted preliminary field research on local water sources, coordinated access to the USGS spatio-temporal dataset, and performed initial exploratory data analysis. Her domain knowledge in environmental science proved valuable for interpreting results.

Shourav Awasthi developed visualization code, formatted the final report, and prepared presentation materials. He ensured reproducibility by maintaining the GitHub repository structure and created all figures used in the report and presentation.

B. Collaboration Challenges and Solutions

The team encountered two significant collaboration challenges:

1) Coordinating asynchronous work: With varying schedules, maintaining synchronization proved difficult initially. We addressed this by establishing a shared Google Colab environment with automated GitHub commits, enabling transparent version control and reducing merge conflicts.

2) Balancing technical depth with accessibility: Ensuring all team members understood technical decisions required regular knowledge-sharing sessions. We instituted weekly review meetings where modeling decisions were explained and discussed, strengthening collective understanding and identifying potential improvements.

VII. IMPACT AND USE OF AI TOOLS

A. AI Tool Usage Declaration

In accordance with academic integrity requirements, we declare the following use of AI-assisted tools:

ChatGPT (OpenAI GPT-4.5): Used responsibly for:

- Report formatting, grammar refinement, and LaTeX/Markdown structure suggestions
- Code documentation and inline comment generation for complex functions
- Literature search assistance for identifying relevant water quality research papers
- Debugging assistance for specific Python errors encountered during implementation

Critical limitation: All experimental design, model training, hyperparameter selection, threshold optimization, and performance analysis were conducted independently by the team. AI tools were not used for generating core analytical results or interpreting findings.

B. Impact Assessment

AI tools provided substantial efficiency gains, particularly in:

- **Documentation quality:** Grammar checking and structural suggestions improved report clarity and professionalism.
- **Code readability:** Automated comment generation helped maintain well-documented code for reproducibility.
- **Debugging efficiency:** Rapid identification of syntax errors and logic bugs accelerated development cycles.

However, AI assistance also introduced risks:

- **Over-reliance:** Initial tendency to accept AI suggestions without verification required conscious effort to maintain critical evaluation.
- **Homogenization:** AI-generated text sometimes lacked the specificity needed for technical reporting, requiring substantial human editing.

VIII. LIMITATIONS AND FUTURE WORK

A. Current Limitations

Several limitations constrain the generalizability and applicability of this work:

1) Dataset constraints: The classification dataset lacks temporal information and geographic metadata, preventing analysis of seasonal variations or location-specific contamination patterns.

2) Feature completeness: Important water quality parameters such as heavy metals, pesticides, and bacterial contamination are absent from the dataset, limiting real-world applicability.

3) Threshold generalization: The optimized threshold ($T = 0.36$) was derived from a single test set and may not generalize to new geographic regions or seasonal conditions without recalibration.

B. Proposed Extensions

Future research should address these limitations through:

1) Deep learning approaches: Recurrent neural networks (LSTMs) or Transformer architectures could capture more complex temporal dependencies in the pH forecasting task, potentially improving beyond the current $R^2 = 0.8329$.

2) Cost-sensitive learning: Explicit incorporation of misclassification costs (quantifying the relative harm of false positives vs. false negatives) could enable more principled threshold selection than F1-score maximization alone.

3) Real-world deployment: Integration with low-cost IoT sensors and edge computing devices for real-time monitoring in Kurukshetra water sources, validating model performance on truly independent data.

4) Uncertainty quantification: Bayesian ensemble methods or conformal prediction could provide confidence intervals around predictions, enabling risk-aware decision making.

IX. CONCLUSION

This research demonstrates that machine learning can substantially enhance water quality prediction in resource-constrained settings. Through systematic threshold optimization, we achieved a 28.6% improvement in F1-score over the baseline, with recall reaching 90% for detecting potable water. These results validate the hypothesis that ensemble methods, properly calibrated, can overcome the challenges of imbalanced environmental classification.

The strong pH forecasting performance ($R^2 = 0.8329$) further demonstrates the feasibility of predictive monitoring systems. Together, these findings provide a reproducible framework for intelligent water quality assessment applicable to similar contexts globally.

Most importantly, the interpretable nature of the Random Forest model—with clear feature importance rankings aligned with water quality science—positions this work for practical deployment in public health decision-making. The high recall achieved ensures that contaminated water samples are rarely missed, prioritizing human safety above all else.

ACKNOWLEDGMENT

The authors thank the Department of Computer Applications at NIT Kurukshetra for computational resources and project guidance. We acknowledge Aditya Kadiwal for making the Water Potability Dataset publicly available on Kaggle, and the USGS for maintaining comprehensive water quality monitoring systems.

REFERENCES

- [1] World Health Organization, *Guidelines for Drinking-Water Quality*, 4th Ed., Geneva, Switzerland, 2017.
- [2] A. Ahmed et al., "Machine learning methods for better water quality prediction," *Journal of Hydrology*, vol. 578, p. 124084, 2019.
- [3] P. Chen, L. Li, and H. Zhang, "Water quality monitoring using machine learning: A comprehensive review," *IEEE Access*, vol. 8, pp. 175020-175040, 2020.
- [4] A. Kadiwal, "Water Potability Dataset," Kaggle, 2020. [Online]. Available: <https://www.kaggle.com/datasets/adityakadiwal/water-potability>
- [5] L. Zhao, O. Gkountouna, and D. Pfoser, "Spatial auto-regressive dependency interpretable learning based on spatial topological constraints," *ACM Transactions on Spatial Algorithms and Systems*, vol. 5, no. 3, pp. 1-28, 2019.
- [6] D. W. Hosmer Jr., S. Lemeshow, and R. X. Sturdivant, *Applied Logistic Regression*, 3rd ed. Hoboken, NJ: Wiley, 2013.
- [7] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [8] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, CA, 2016, pp. 785-794.
- [9] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowledge and Data Engineering*, vol. 21, no. 9, pp. 1263-1284, Sep. 2009.
- [10] G. Louppe, "Understanding random forests: From theory to practice," Ph.D. dissertation, Dept. Elect. Eng. Comput. Sci., Univ. Liège, Belgium, 2014.
- [11] S. S. Baral, S. N. Das, and P. Rath, "Hardness removal from drinking water using low-cost adsorbents," *Journal of Environmental Science and Engineering*, vol. 49, no. 4, pp. 253-258, 2007.
- [12] M. M. Rahman and M. H. Kabir, "Water quality assessment and classification of Dhaka's groundwater: A machine learning approach," *Environmental Monitoring and Assessment*, vol. 192, p. 741, 2020.