

Potable Water Identification with Machine Learning: An Exploration of Water Quality Parameters

¹B R Mohan, ²Dr. Dileep M, ³Dr Vijay Bhuria, ⁴Sai Sudha Gadde, ⁵Dr. Kumarasamy M, ⁶Achyutha Prasad N

¹Computer Science and Engineering, East West Institute of Technology, Bangalore
brmohan398@gmail.com

²Strategy, Faculty of Management Sciences, Nile University of Nigeria, Abuja.
Email: Dileep.KM@nileuniversity.edu.ng

³Electrical Department, Madhav Institute of Technology & Science, Gwalior, India.
Email: vijay.bhuria@mitsgwalior.in

⁴Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur District, A.P., India.
Email: Saisudhagadde@gmail.com

⁵Department of Computer Science, College of Engineering and Technology,
Wollega University, Nekemte, Oromia Region, Ethiopia.
Email: drmsamy115@gmail.com

⁶Department of Computer Science and Engineering, East West Institute of Technology, Bangalore, India
Email: achyuth001@gmail.com

Abstract- In this research, we aim to determine the water potability using three machine learning classification algorithms: decision tree, gradient boosting and bagging classifier. These algorithms were trained and tested on a dataset of water quality measurements. The outcomes of the experiment showed that the gradient boosting algorithm achieved the highest F1-score of 0.78 among all the algorithms. This indicates that the gradient boosting algorithm was most effective in correctly identifying both the safe and contaminated water samples. The results of this study demonstrate that gradient boosting is a promising approach for determining water potability and can be used as a reliable method for water quality assessment.

Keywords- Water Quality; Machine Learning; Water Potability; Classification

I. INTRODUCTION

Water potability refers to the quality of water that is safe for human consumption. It is a crucial aspect of public health and sanitation, such as having access to sanitary drinking water is essential for maintaining a healthy lifestyle. Water contamination can occur due to a variety of factors such as industrial and agricultural pollution, natural disasters, and improper sewage treatment. Contaminated water can cause a range of health issues, including diarrhea, cholera, and typhoid fever [1-4]. In severe cases, it can lead to death. To ensure water potability, water is tested for a variety of contaminants, including bacteria, viruses, and chemicals such as lead and arsenic. These tests are conducted by government agencies and independent organizations, and the results are made available to the public. However, even with regular testing, it can be difficult to detect all contaminants in water [5-8]. This is where machine learning comes in. Machine learning is able to recognise patterns and trends in water quality data that may not be immediately obvious to human researchers by applying sophisticated algorithms. In the end, enhanced water potability can result from the more precise and effective detection of water pollutants. In addition to the

public health benefits, ensuring water potability also has economic benefits. Because it supports increased agricultural production, industrial development, and tourism, access to clean, safe drinking water is crucial for economic progress. Overall, water potability is a critical issue that affects the health and well-being of individuals and communities. The use of machine learning in potability detection can help improve the accuracy and efficiency of water testing, ensuring access to clean and safe drinking water for all. Machine learning is a rapidly growing field that is revolutionizing many industries, including manufacturing, healthcare, and transportation. In manufacturing, machine learning is used to improve productivity and efficiency [9-15]. For example, machine learning algorithms can be utilized to optimize production schedules, forecast equipment failures, and improve product quality. Additionally, machine learning can be used to improve supply chain management by predicting demand and optimizing inventory levels. These benefits are resulting in cost savings and increased competitiveness for manufacturing companies. [16-20] In healthcare, machine learning is used to improve patient outcomes and reduce costs. For example, machine learning

algorithms can be used to envisage patient outcomes, identify patients at high risk of complications, and optimize treatment plans. Additionally, machine learning can be used to improve medical imaging by identifying diseases and detecting abnormalities [21-24]. As a result, healthcare providers are able to deliver more effective and efficient care to patients. In transportation, machine learning is used to improve safety and efficiency. For instance, route optimization and traffic pattern prediction using machine learning algorithms can improve travel times and reduce congestion [25-28]. Additionally, machine learning can be utilized to advance the protection of autonomous vehicles by detecting and responding to potential hazards. These benefits are helping to make transportation more efficient and safer for everyone.

II. PROBLEM STATEMENT

The availability of clean and safe drinking water is a basic human need, yet in many places of the world, getting access to potable water is still very difficult. Traditional methods for detecting water potability, such as chemical testing and visual inspections, can be time-consuming and costly, and may not provide accurate results in all cases.

The goal of this research is to create a machine learning model that can quickly and reliably determine if water is potable. The model should be able to analyze data from a variety of sources, including water samples, sensor readings, and environmental factors, to identify patterns and predict the presence of contaminants.

The proposed solution will leverage the power of machine learning algorithms to detect water potability, with the goal of providing a more accurate and efficient method for ensuring the safety of drinking water. Additionally, the model will be able to learn from data over time, allowing it to continuously improve its predictions and adapt to changing water conditions.

This work will have a significant impact on the public health and well-being, by providing a powerful tool for detecting water potability and ensuring access to clean and safe drinking water for communities around the world. It will also help to reduce the costs associated with traditional methods of water potability testing, and make it more accessible for developing countries.

III. EXPERIMENTAL PROCEDURE

A key element of effective health protection measures and a fundamental human right, access to clean drinking water is essential for good health. This is important as a health and development issue on a national, regional, and local level. Because they lower unfavourable health consequences and medical costs more than they cost to implement, investments in water supply and sanitation have been demonstrated to generate a net economic advantage in some locations. In the present work, the dataset file includes measurements of water quality for 3276 dissimilar frames of water. The input parameters are pH, Hardness, Solids (Total dissolved solids - TDS), Chloramines, Sulfate, Conductivity, Organic carbon, Trihalomethanes, Turbidity while Potability will be the output parameter. The given dataset will be subjected to classification-based machine learning algorithms for the forecast of the potability of the water. In the present work, the output parameter i.e. Water potability is labelled as 0 if its not potable and labelled as 1 if the water is potable. Supervised machine learning classification algorithms are a type of algorithm used to classify data into different categories or classes. These algorithms work by training a model on a labeled dataset, where the data points are already associated with a specific class or label. The algorithm's objective is to identify the underlying links and patterns in the data so that it may utilise this understanding to anticipate the outcomes of new, upcoming data. The process of training a supervised machine learning classification algorithm typically begins with the selection of a suitable algorithm, such as decision trees, logistic regression, or support vector machines. Then, a training set and a testing set are created from the labelled dataset. The model is trained using the training set, and its performance is assessed using the testing set. During the training process, the algorithm iteratively analyzes the data and learns the patterns and relationships that are associated with each class. This is done by adjusting the model's parameters, such as the weights of the features, to minimize the error between the predicted and true labels. Once the training is complete, the model is tested using the testing set, and the performance is appraised using system of measurement such as accuracy, precision, recall, and F1-score. If the performance is satisfactory, the model can then be used to classify new, unseen data points. The implemented framework used in the present work is shown in Figure 1.

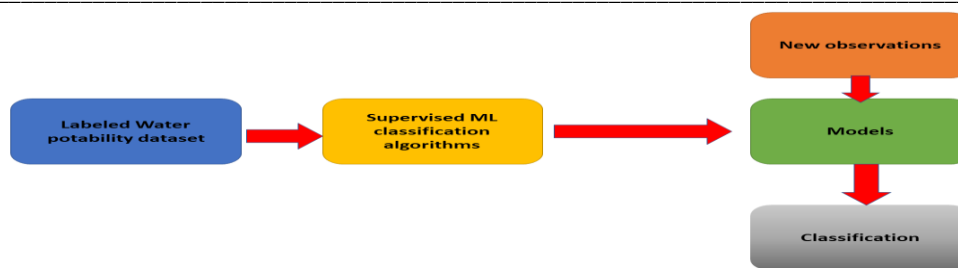


Figure 1. Machine Learning classification-based framework implemented in the present work

RESULTS AND DISCUSSION

Table 1 shows the statistical analysis results of the present dataset.

TABLE I. STATISTICAL ANALYSIS RESULTS

| a~ | mean | std | min | 25% | 50% | 75% | max | |
|-----------------|-------------|--------------|-------------|------------|--------------|--------------|--------------|--------------|
| ph | 2785.000000 | 7.080795 | 1.594320 | 0.000000 | 6.093092 | 7.036752 | 8.062066 | 14.000000 |
| Hardness | 3276.000000 | 196.369496 | 32.879761 | 47.432000 | 176.850538 | 196.967627 | 216.667456 | 323.124000 |
| Solids | 3276.000000 | 22014.092526 | 8768.570828 | 320.942611 | 15666.690297 | 20927.833607 | 27332.762127 | 61227.196008 |
| Chloramines | 3276.000000 | 7.122277 | 1.583085 | 0.352000 | 6.127421 | 7.130299 | 8.114887 | 13.127000 |
| Sulfate | 2495.000000 | 333.775777 | 41.416840 | 129.000000 | 307.699498 | 333.073546 | 359.950170 | 481.030642 |
| Conductivity | 3276.000000 | 426.205111 | 80.824064 | 181.483754 | 365.734414 | 421.884968 | 481.792304 | 753.342620 |
| Organic_carbon | 3276.000000 | 14.284970 | 3.308162 | 2.200000 | 12.065801 | 14.218338 | 16.557652 | 28.300000 |
| Trihalomethanes | 3114.000000 | 66.396293 | 16.175008 | 0.738000 | 55.844536 | 66.622485 | 77.337473 | 124.000000 |
| Turbidity | 3276.000000 | 3.966786 | 0.780382 | 1.450000 | 3.439711 | 3.955028 | 4.500320 | 6.739000 |

Statistical analysis is a key component of machine learning, as it is used to appreciate, interpret, and type predictions grounded on the data. By applying statistical methods, machine learning algorithms can extract useful information and patterns from data, and use this knowledge to make predictions or decisions.

One of the main results obtained by statistical analysis in machine learning is the estimation of model parameters. This is done by fitting a statistical model to the data, which involves finding the values of the parameters that best explain the data. Finding coefficient values that minimise the sum of squared errors between the predicted and true values, for instance, is the objective in linear regression. Another important result obtained by statistical analysis in machine learning is the assessment of model performance. This is accomplished by assessing the model's predictions' accuracy using various measures, including accuracy, precision, recall, and F1-score. These metrics give a mechanism to measure how well the model can categorise or forecast the results of fresh, unforeseen data points.

Statistical analysis also helps in understanding the feature importance and variable selection. By performing statistical tests, like chi-square test, we can understand which feature has a stronger relationship with the target variable. This helps us in selection of important features to be considered while building the model and also helps in reducing the dimensionality of the data.

Exploratory Data Analysis (EDA), which aids in understanding the underlying patterns and relationships within the data, is a vital phase in the machine learning process. By using EDA, we may find any outliers, missing numbers, and other anomalies in the data as well as learn more about the distribution and structure of the data at its core. Finding significant features is one of the key outcomes of EDA in machine learning. By visualizing the data and calculating summary statistics, we can gain insights into which features are most important for the problem at hand. This information can then be used to guide feature selection and dimensionality reduction. Another key result obtained by EDA in machine learning is the understanding of the relationship between features and the target variable. By

creating visualizations such as scatter plots and heatmaps, we can identify any linear or non-linear relationships between the features and the target variable. These insights can help in selecting the appropriate model and fine-tuning the model's parameters.

EDA also helps in understanding the distribution of data and identifying any outliers. By visualizing the data distribution,

we can identify any data points that fall outside the range of normal values. This is particularly useful for identifying errors in the data and for making decisions about data cleaning. Figure 2 shows the feature distribution plot obtained in the current work. Figure 3 shows the obtained pair plot in the current work.

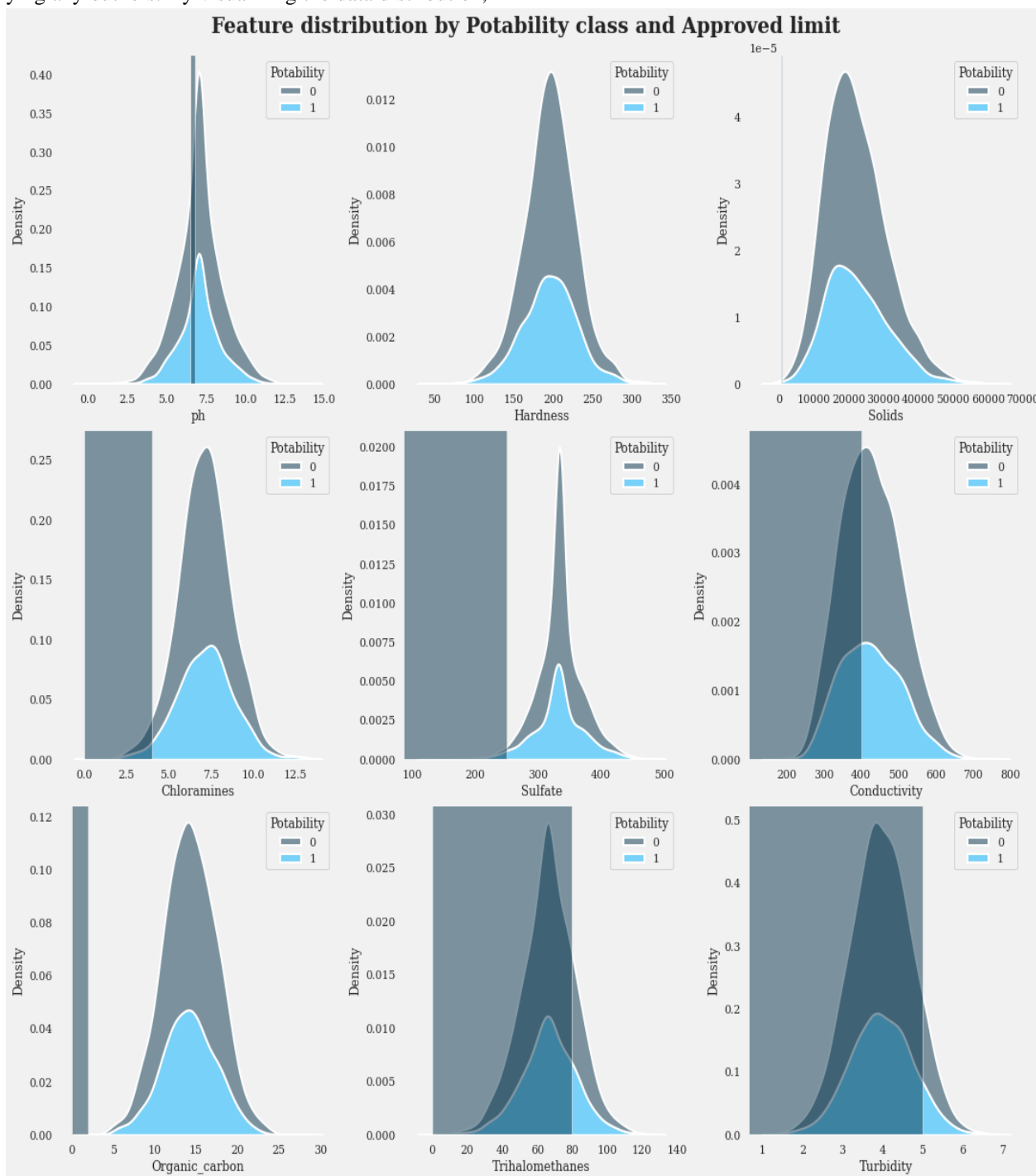


Figure 2. Feature distribution plot obtained in the current work

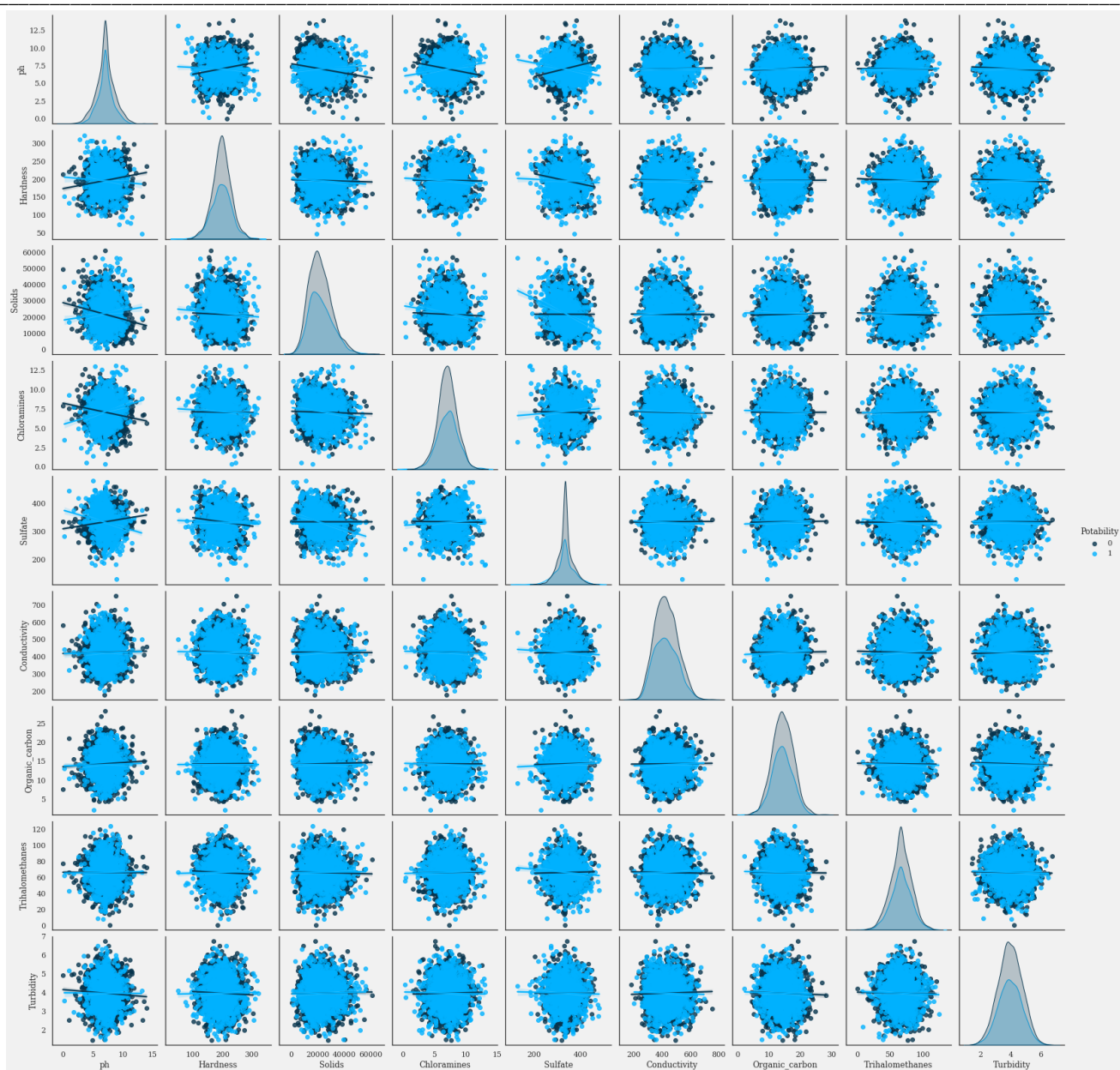


Figure 3. Obtained pair plot in the present work

For the purpose of determining the potability of water in the current work, three machine learning classification-based algorithms—Decision Tree, Gradient Boosting, and Bagging Classifier—have been applied. A popular supervised learning technique for classification and regression applications is the decision tree. It is a decision-making model that resembles a tree with potential outcomes. Recursively dividing the dataset into subgroups according to the values of an input feature is how the tree is constructed. A test on an input feature is represented by each internal node of the tree, the test's result is represented by each branch, and the class label is represented by each leaf node. The characteristic that best divides the dataset into subsets is chosen as the root node to start the decision tree classification process. Using a criterion

like information gain or gain ratio, the features are chosen. The root node is chosen as the characteristic having the highest value of the criterion.

The feature that divides each subset most effectively is then chosen to further divide the subsets produced by the root node. As long as there are samples of the same class in all of the subsets, this process is repeated recursively for each internal node. The class labels of the samples in the respective subsets are represented by the leaf nodes. By moving up the tree from its root to one of its leaf nodes, the tree can be utilised for categorization. The anticipated class label for a new sample is the class label of the leaf node. Decision trees have the benefit of being simple to comprehend and interpret,

which is one of its key advantages. The model's decisions are clearly and understandably explained by the tree structure. As a result, decision trees are frequently used in combination with ensemble methods like random forest, which use several decision trees to arrive at a final conclusion. Decision trees, however, can be sensitive to slight alterations in the data and can overfit readily.

An ensemble machine learning method called gradient boosting is employed for both classification and regression problems. The method is built on the notion of creating a model iteratively by including weak learners who fix the errors of the prior models. Decision trees are frequently used in gradient boosting as weak learners. Gradient boosting classification's mechanism begins with a straightforward model, like a decision tree with a single split. On the basis of the training data, this model is used to create predictions, and the errors are determined. The following stage is to match a fresh model to the flaws in the old one. Iterations of this process are repeated indefinitely, or until the model's performance plateaus. The new model concentrates on the data that the old model misclassified in each iteration. This is accomplished by giving the incorrectly identified samples a greater weight. The new model is taught to fix the errors committed by the old model. The final model is the weighted sum of each weak model, with each model's weight based on how well it performed on the training set of data. When compared to a single decision tree, gradient boosting exhibits state-of-the-art performance on a variety of classification and regression problems and is robust to overfitting. This is so that the variance can be reduced by the final model, which combines numerous weak models. However, it can be

Table II. shows the obtained results for measuring the performance of the implemented algorithms.

Table II. OBTAINED RESULTS OF THE ML ALGORITHMS

| Algorithms | Precision for 0 and 1 | Recall for 0 and 1 | Overall F1-Score |
|---------------------------|-----------------------|--------------------|------------------|
| Decision Tree | 0.66 and 0.56 | 0.90 and 0.22 | 0.65 |
| Gradient Boosting | 0.86 and 0.67 | 0.77 and 0.80 | 0.78 |
| Bagging Classifier | 0.81 and 0.68 | 0.80 and 0.69 | 0.76 |

It is observed that Gradient boosting is resulting in the highest F1-Score value. Gradient boosting is a powerful ensemble method that combines multiple weak learners to create a strong classifier. It is a popular choice for classification problems because of its ability to achieve high accuracy. The main reason why gradient boosting is so effective for classification is its ability to adapt to non-linear decision boundaries. Gradient boosting uses decision trees as its base learners, which are known for their ability to model complex and non-linear relationships in the data. By combining

sensitive to the selection of hyperparameters and is also computationally expensive.

A bagging classifier is an ensemble technique that combines a number of base classifiers in order to enhance the model's overall performance and lower its variance. The basis classifiers are trained using several bags created from the training data. Through majority vote or by averaging the predictions of the basic classifiers, the final prediction is determined. The mechanism of bagging classifier begins with selecting a random subset of the training data, with replacement, to train each base classifier. This is done to introduce randomness and diversity among the base classifiers. The original dataset's data points are randomly chosen using replacement to form the subsets, sometimes referred to as bags. This implies that while certain data points might be chosen more than once in a bag, others might not be chosen at all. The basis classifiers are trained on these subsets of data after the bags have been produced. Combining the basis classifiers' predictions yields the final prediction. There are two ways to combine the predictions: majority voting and averaging. In majority voting, the class that is predicted by the majority of the base classifiers is chosen as the final prediction. In averaging, the predictions of the base classifiers are averaged to produce the final prediction.

Bagging classifiers are known to be effective in reducing the variance of the model. By training the base classifiers on different subsets of data, the model is able to capture different patterns in the data, which helps to reduce the variance of the model. Bagging classifiers are also known to be robust to noise in the data.

multiple decision trees, gradient boosting is able to capture more intricate patterns in the data, resulting in a more accurate classifier. Another key aspect of gradient boosting that contributes to its high accuracy is its ability to identify and correct errors made by previous base learners. The method does this by focusing on the samples that were misclassified by previous base learners and giving more weight to those samples when training the next base learner. This allows the classifier to improve upon its performance with each iteration. Gradient boosting also has a built-in feature

selection mechanism which is known as "shrinkage", which helps to reduce overfitting. Shrinkage reduces the influence of each base learner by adjusting the learning rate of the model, which helps to stabilize the final predictions. This is particularly useful when working with high-dimensional data, as it helps to reduce the risk of overfitting.

IV. CONCLUSION

In summary, this study used three machine learning classification algorithms—the decision tree, gradient boosting, and bagging classifier—to assess the potability of the water. The results of the experiment showed that the gradient boosting algorithm achieved the uppermost F1-score of 0.78 among all the algorithms, indicating that it was most effective in correctly identifying both the safe and contaminated water samples. These results demonstrate that gradient boosting is a promising approach for determining water potability and can be used as a reliable method for water quality assessment. It is crucial to remember that this study was based on a particular dataset, and that in order to validate the results, additional research should be done using different datasets and under different circumstances.

As for future work, one could explore using other ensemble methods like Random Forest or LightGBM and compare the performance with Gradient Boosting. Additionally, incorporating other features like location, weather conditions etc. that may have an impact on water quality can be added to the dataset to evaluate if the model's accuracy increases. As an alternative to conventional machine learning techniques, the usage of deep learning models like Convolutional Neural Networks (CNNs) can also be investigated.

REFERENCES

- [1] Sharma, S. and Bhattacharya, A.J.A.W.S., 2017. Drinking water contamination and treatment techniques. *Applied water science*, 7(3), pp.1043-1067.
- [2] Pye, V.I. and Patrick, R., 1983. Ground water contamination in the United States. *Science*, 221(4612), pp.713-718.
- [3] Schweitzer, L. and Noblet, J., 2018. Water contamination and pollution. In *Green chemistry* (pp. 261-290). Elsevier.
- [4] Bedient, P.B., Rifai, H.S. and Newell, C.J., 1994. *Ground water contamination: transport and remediation*. Prentice-Hall International, Inc..
- [5] Gong, J., Guo, X., Yan, X. and Hu, C., 2023. Review of Urban Drinking Water Contamination Source Identification Methods. *Energies*, 16(2), p.705.
- [6] Vanlalmingmawia, C., Lee, S.M. and Tiwari, D., 2023. Plasmonic noble metal doped titanium dioxide nanocomposites: newer and exciting materials in the remediation of water contaminated with micropollutants. *Journal of Water Process Engineering*, 51, p.103360.
- [7] Qin, H. and Doll, G.L., 2023. Effects of Water Contamination on Micropitting and Rolling Contact Fatigue of Bearing Steels. *Journal of Tribology*, 145(1), p.011501.
- [8] Bradley, P.M., Romanok, K.M., Smalling, K.L., Focazio, M.J., Evans, N., Fitzpatrick, S.C., Givens, C.E., Gordon, S.E., Gray, J.L., Green, E.M. and Griffin, D.W., 2023. Bottled water contaminant exposures and potential human effects. *Environment International*, 171, p.107701.
- [9] Mishra, A., 2020. Artificial intelligence algorithms for the analysis of mechanical property of friction stir welded joints by using python programming. *Welding Technology Review*, 92.
- [10] Wuest, T., Weimer, D., Irgens, C. and Thoben, K.D., 2016. Machine learning in manufacturing: advantages, challenges, and applications. *Production & Manufacturing Research*, 4(1), pp.23-45.
- [11] Rai, R., Tiwari, M.K., Ivanov, D. and Dolgui, A., 2021. Machine learning in manufacturing and industry 4.0 applications. *International Journal of Production Research*, 59(16), pp.4773-4778.
- [12] Thapliyal, S. and Mishra, A., 2021. Machine learning classification-based approach for mechanical properties of friction stir welding of copper. *Manufacturing Letters*, 29, pp.52-55.
- [13] Pham, D.T. and Afify, A.A., 2005. Machine-learning techniques and their applications in manufacturing. *Proceedings of the Institution of Mechanical Engineers, Part B: Journal of Engineering Manufacture*, 219(5), pp.395-412.
- [14] Wang, S., Shen, Z., Shen, Z., Dong, Y., Li, Y., Cao, Y., Zhang, Y., Guo, S., Shuai, J., Yang, Y. and Lin, C., 2021. Machine-learning micropattern manufacturing. *Nano Today*, 38, p.101152.
- [15] Jatti, V.S., Dhabale, R.B., Mishra, A., Khedkar, N.K., Jatti, V.S. and Jatti, A.V., 2022. Machine Learning Based Predictive Modeling of Electrical Discharge Machining of Cryo-Treated NiTi, NiCu and BeCu Alloys. *Applied System Innovation*, 5(6), p.107.
- [16] Hu, H., Xu, J., Liu, M. and Lim, M.K., 2023. Vaccine supply chain management: An intelligent system utilizing blockchain, IoT and machine learning. *Journal of Business Research*, 156, p.113480.
- [17] Rosenberg-Vitorica, W., Salais-Fierro, T.E., Marmolejo-Saucedo, J.A. and Rodriguez-Aguilar, R., 2023. Machine Learning Applications in the Supply Chain, a Literature Review. *Smart Applications with Advanced Machine Learning and Human-Centred Problem Design*, pp.753-761.
- [18] Garinian, E.O.M., Fierro, T.E.S., Saucedo, J.A.M. and Aguilar, R.R., 2023. Machine Learning Applications for Demand Driven in Supply Chain: Literature Review. *Smart Applications with Advanced Machine Learning and Human-Centred Problem Design*, pp.763-772.
- [19] Hamdan, I.K., Aziguli, W., Zhang, D. and Sumarlah, E., 2023. Machine learning in supply chain: prediction of real-time e-order arrivals using ANFIS. *International*

- Journal of System Assurance Engineering and Management, pp.1-20.
- [20] Panigrahi, R.R., Dash, M., Shaikh, Z.H. and Irfan, M., 2023. Review of Machine Learning Techniques in the Supply Chain Management of Indian Industry: A Future Research Agenda. In *Advanced Machine Learning Algorithms for Complex Financial Applications* (pp. 199-219). IGI Global.
- [21] Mooney, S.D., 2023. Technology Platforms and Approaches for Building and Evaluating Machine Learning Methods in Healthcare. *The Journal of Applied Laboratory Medicine*, 8(1), pp.194-202.
- [22] Roberts, L., Dhanoa, H., Lanes, S. and Holdship, J., 2023. Machine learning for enhanced healthcare: an overview for operational and clinical leads. *British Journal of Healthcare Management*, 29(1), pp.12-19.
- [23] Zini, M. and Carcasci, C., 2023. Machine learning-based monitoring method for the electricity consumption of a healthcare facility in Italy. *Energy*, 262, p.125576.
- [24] Rayan, R.A., 2023. Machine Learning for Smart Health Care. *Machine Learning Algorithms and Applications in Engineering*, p.1.
- [25] Megnidio-Tchoukouegno, M. and Adedeji, J.A., 2023. Machine Learning for Road Traffic Accident Improvement and Environmental Resource Management in the Transportation Sector. *Sustainability*, 15(3), p.2014.
- [26] Rolczynski, B.S., Díaz, S.A., Kim, Y., Mathur, D., Klein, W.P., Medintz, I. and Melinger, J., 2023. Determining Interchromophore Effects for Energy Transport in Molecular Networks Using Machine-Learning Algorithms. *Physical Chemistry Chemical Physics*.
- [27] Vadlamani, S. and Modashiya, M., 2023. Improving Transportation Planning Using Machine Learning. In *Encyclopedia of Data Science and Machine Learning* (pp. 3076-3088). IGI Global.
- [28] Jayalakshmi, S., SV, A.K., Ranga, J., Venkatasubramanian, K. and Lavanya, Y., 2023. A Machine Learning-Based Image Segmentation for Real-Time Images in Smart Intelligent Transportation Systems. In *Handbook of Research on Advanced Practical Approaches to Deepfake Detection and Applications* (pp. 147-161). IGI Global.
- [29] Achyutha Prasad, N., Guruprakash, C.D., 2019. A relay mote wheeze for energy saving and network longevity enhancement in WSN. *International Journal of Recent Technology and Engineering* 8, 8220–8227. doi:10.35940/ijrte.C6707.098319.
- [30] Kalshetty, J. N., Achyutha Prasad, N., Mirani, D., Kumar, H., & Dhingra, H. (2022). Heart health prediction using web application. *International Journal of Health Sciences*, 6(S2), 5571–5578. <https://doi.org/10.53730/ijhs.v6nS2.6479>.
- [31] Achyutha, P. N., Hebbale, S., & Vani, V. (2022). Real time COVID-19 facemask detection using deep learning. *International Journal of Health Sciences*, 6(S4), 1446–1462. <https://doi.org/10.53730/ijhs.v6nS4.6231>