

Machine Learning Model for Water Quality Prediction using Python and AI framework

Dr. Kalaivazhi Vijayaragavan¹, N. Praveen², M. V. Sudharsan³ and P. S. Vijayan⁴

Associate Professor, Department of Information Technology¹

Students, B.Tech., Final Year, Department of Information Technology^{2,3,4}

Anjalai Ammal Mahalingam Engineering College, Thiruvavur, India

Abstract: During the last years, water quality has been threatened due to unprocessed effluents, municipal refuse, factory wastes, junking of compostable and non-compostable effluents has hugely contaminated nature-provided water bodies like rivers, lakes and ponds are pollutants. Therefore, it is necessary to look into the water standards before the usage. Hence modeling and predicting water quality have become very important in controlling water pollution. Safe drinking-water access is essential to health, a basic human right and a component of effective policy for health protection. It is important as a health and development issue at a national, regional and local level. Thus it is a problem that can greatly benefit from Artificial Intelligence (AI). Traditional methods require human inspection and is time consuming. Automatic Machine Learning (AutoML) facilities provide machine learning with push of a button, or, on a minimum level, ensure to retain algorithm execution, data pipelines, and code, generally, are kept from sight and are anticipated to be the stepping stone for normalizing AI. However, it is a field under research still. This project work aims to recognize the areas where an AutoML system falls short or outperforms a traditional expert system built by data scientists. Keeping this as the motive, this project work dives into the Machine Learning (ML) algorithms for comparing AutoML and an expert architecture built by this project for Water Quality Assessment to evaluate the Water Quality Index, which gives the general water quality, and the Water Quality Class, a term classified on the basis of the Water Quality Index using python. In this Project, we are going to implement a water quality prediction using machine learning techniques. In this project, our model predicts, that the water is safe to drink or not, using some parameters like PH value, conductivity, hardness, etc. Finally the results of accuracy level of AutoML and Python compared with conventional ML techniques.

Keywords: Machine Learning, Classification Algorithm, Prediction, PyThon and AI framework

I. INTRODUCTION

Machine learning is an application of AI that enables systems to learn and improve from experience without being explicitly programmed. Machine learning focuses on developing computer programs that can access data and use it to learn for themselves. Similar to how the human brain gains knowledge and understanding, machine learning relies on input, such as training data or knowledge graphs, to understand entities, domains and the connections between them. With entities defined, deep learning can begin. The machine learning process begins with observations or data, such as examples, direct experience or instruction. It looks for patterns in data so it can later make inferences based on the examples provided. The primary aim of ML is to allow computers to learn autonomously without human intervention or assistance and adjust actions accordingly. Machine learning as a concept has been around for quite some time. The term “machine learning” was coined by Arthur Samuel, a computer scientist at IBM and a pioneer in AI and computer gaming. Samuel designed a computer program for playing checkers. The more the program played, the more it learned from experience, using algorithms to make predictions.

II. PYTHON AND AI FRAMEWORK

- **Python:** Python is a computer programming language often used to build websites and software, automate tasks, and conduct data analysis. Python is a general-purpose language, meaning it can be used to create a variety of different programs and isn't specialized for any specific problems. This versatility, along with its beginner-

friendliness, has made it one of the most-used programming languages. A survey conducted by industry analyst found that it was the second-most popular programming language among developers in 2021.

- **AutoML:** Automated machine learning is the process of applying machine learning models to use real-world problems using automation. More specifically, it automates the selection, composition and parameterization of Machine Learning models. Automating the ML process makes it more user-friendly and often provides faster, more accurate outputs than hand-coded algorithms. AutoML is a typically platform or open source library that simplifies each step in the ML process, from handling a raw dataset to deploying a practical ML model. In traditional ML, models are developed by hand, and each step in the process must be handled separately.

III. DESIGN ISSUES

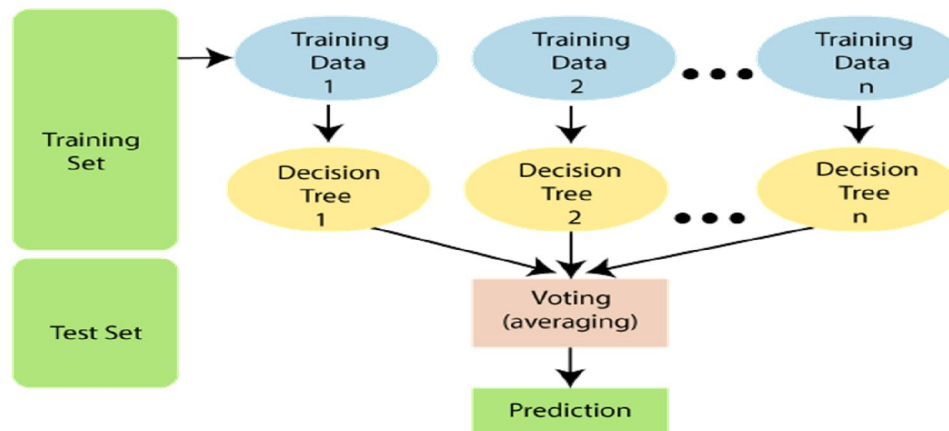
The challenge is aimed to make use of machine learning algorithm in Water Quality Assessment to evaluate the Water Quality Index of the dataset. In this project, we aim to impart the ability to get rid of biases in a machine algorithm and to predict the accuracy of the datasets.

- To evaluate the training speed of AutoML and Python based on Classification Algorithm.
- Design of a machine learning model, which can classify the different datasets.
- Datasets using Supervised and Unsupervised Learning techniques analyses the accuracy of the water quality based on parameters like PH value, conductivity and hardness.
- Machine learning algorithms use different methods to analyse training data and apply what they learn to new examples.
- When choosing a machine learning framework, it is important to consider whether this adjustment should be automatic or manual.
- AutoML library and Python platform to work with deep neural networks, testing array operations in order to get better accuracy.

3.1 Algorithm Used

A. Random Forest Classification:

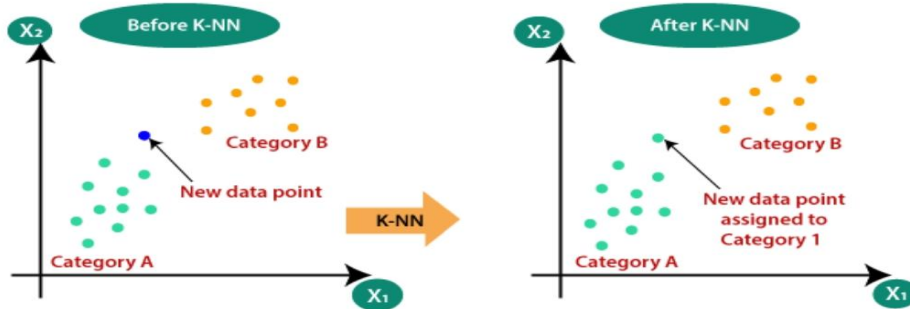
- A random forest is a machine learning technique, that is used to solve regression and classification problems. It utilizes ensemble learning, which is technique that combines many classifiers to provide solutions to complex problems.
- A random forest Classification algorithm consists of many decision trees. The ‘forest’ generated by the random forest classification algorithm is trained through bagging or bootstrap aggregating. Bagging is an ensemble meta-algorithm that improves the accuracy of ML algorithms.
- The (random forest classification) algorithm establishes the outcome based on the predictions of the decision trees. It predicts by taking average or mean of the output from various trees. Increasing the number of trees and increases the precision of the outcome.





B. K-Nearest Neighbour:

K Nearest Neighbor algorithm(KNN) falls under the Supervised Learning category and is used for classification and regression. It is a versatile algorithm and used for imputing missing values also resampling datasets. As the name (KNN) suggests it considers K Nearest Neighbors to predict the class or continuous value for the new Datapoint.



C. In AutoML using Tpot

TPOT (Tree-based Pipeline Optimization Tool) is a AutoML tool specifically designed for the efficient construction of optimal pipelines through genetic programming. TPOT is a open source library and makes use of scikit-learn components for data transformation, feature decomposition, feature selection and model selection .Although TPOT is classified as AutoML tool, as such it does not offer the “end-to-end” of an Machine Learning pipeline. TPOT is merely focused on the optimized automation of specific components of an Machine Learning pipeline. we can see the phases automated by TPOT and the ones specifically addressed by the Data Scientist or Machine Learning Engineer.

3.2 Development Model

The first stage of development of Artificial Intelligence models is the preparation of the dataset. In this stage, the collected dataset shall be divided into two groups, training and testing. The training and testing dataset are used to the calibration and validation of applied models, respectively. Depending on the simulation conditions regarding time series modeling or function fitting, the approach of assigning a dataset for each group are different. In time series modeling, the history of collecting data shall be considered and shuffling the dataset is not correct, whereas for function fitting using data shuffling idea is allowed. Usually for both scenarios, about 70%–80% of the dataset is assigned for calibration and the remaining 20%–30% for validation. The next step for developing the AI models, such as Random forest classification, K-nearest neighbor and Tpot in AutoML is designing the architecture of the network.

3.3 Testing Analysis

We are going to implement a water quality prediction using machine learning techniques. We will implement in this project in Random forest classification and K-nearest neighbor algorithm in supervised learning and Tpot in AutoML. Then we compare python and AI framework, Finally we find which one is accurate the Highest level.

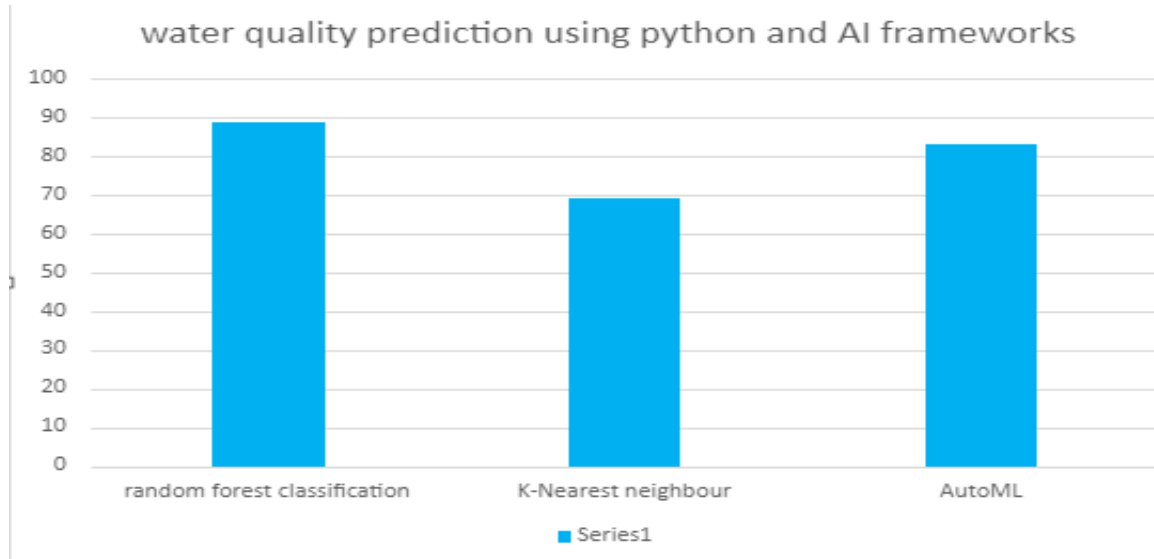
ALGORITHM	ACCURACY LEVEL
Random forest Classification	0.89%
K-nearest neighbor	0.68%
TPOT in AutoML	0.83%



```
Python 3.8.0 Shell
File Edit Shell Debug Options Window Help
Python 3.8.0 (tags/v3.8.0:fa919fd, Oct 14 2019, 19:37:50) [MSC v.1916 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
===== RESTART: C:\Python\wa2.py =====
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3276 entries, 0 to 3275
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 ph 2785 non-null float64
1 Hardness 3276 non-null float64
2 Solids 3276 non-null float64
3 Chloramines 3276 non-null float64
4 Sulfate 2495 non-null float64
5 Conductivity 3276 non-null float64
6 Organic_carbon 3276 non-null float64
7 Trihalomethanes 3114 non-null float64
8 Turbidity 3276 non-null float64
9 Potability 3276 non-null int64
dtypes: float64(9), int64(1)
memory usage: 256.1 KB
Best parameters for KNN: {'n_neighbors': 1}
Best parameters for Random Forest: {'min_samples_leaf': 2, 'n_estimators': 200}
K Nearest Neighbours : 0.68
Random Forest : 0.89
>>> |
```

```
Python 3.8.0 Shell
File Edit Shell Debug Options Window Help
| 156/200 [09:33<04:10, 5.70s/pipeline]Optimization Progress: 78%|██████████| 157/200 [09:35<03:14, 4.52s/pipeline]Optimization
n Progress: 79%|██████████| 158/200 [10:40<15:51, 22.65s/pipeline]Optimization Progress: 80%|██████████| 159/200 [10:42<11:18, 16.5
4s/pipeline]Optimization Progress: 80%|██████████| 160/200 [10:46<08:30, 12.75s/pipeline]Optimization Progress: 80%|██████████| 161
/200 [10:48<06:08, 9.44s/pipeline]Optimization Progress: 81%|██████████| 162/200 [10:53<05:14, 8.26s/pipeline]Optimization Progress:
82%|██████████| 163/200 [11:01<05:04, 8.22s/pipeline]Optimization Progress: 82%|██████████| 164/200 [11:03<03:46, 6.30s/pipeli
ne]Optimization Progress: 82%|██████████| 165/200 [11:11<04:01, 6.90s/pipeline]Optimization Progress: 83%|██████████| 166/200 [11:
13<03:04, 5.43s/pipeline]Optimization Progress: 84%|██████████| 167/200 [11:17<02:37, 4.77s/pipeline]Optimization Progress: 84%|
██████████| 168/200 [11:19<02:05, 3.91s/pipeline]Optimization Progress: 84%|██████████| 169/200 [11:21<01:44, 3.37s/pipeline]Optim
ization Progress: 85%|██████████| 170/200 [11:25<01:47, 3.59s/pipeline]Optimization Progress: 86%|██████████| 171/200 [11:27<01:30
, 3.13s/pipeline]Optimization Progress: 86%|██████████| 172/200 [11:30<01:23, 2.98s/pipeline]Optimization Progress: 86%|██████████|
173/200 [11:31<01:11, 2.66s/pipeline]Optimization Progress: 87%|██████████| 174/200 [12:02<04:47, 11.04s/pipeline]Optimization P
rogress: 88%|██████████| 175/200 [12:06<03:40, 8.83s/pipeline]Optimization Progress: 88%|██████████| 176/200 [12:08<02:42, 6.77s/
pipeline]Optimization Progress: 88%|██████████| 177/200 [12:13<02:23, 6.25s/pipeline]Optimization Progress: 89%|██████████| 178/20
0 [12:15<01:49, 4.96s/pipeline]Optimization Progress: 90%|██████████| 179/200 [12:19<01:40, 4.81s/pipeline]Optimization Progress:
90%|██████████| 180/200 [12:22<01:22, 4.13s/pipeline]Optimization Progress: 90%|██████████| 181/200 [12:36<02:17, 7.25s/pipeline]
Optimization Progress: 91%|██████████| 182/200 [12:41<01:54, 6.38s/pipeline]Optimization Progress: 92%|██████████| 183/200 [12:51<
02:11, 7.76s/pipeline]Optimization Progress: 92%|██████████| 184/200 [12:57<01:55, 7.22s/pipeline]Optimization Progress: 92%|██████████|
185/200 [12:59<01:24, 5.66s/pipeline]Optimization Progress: 93%|██████████| 186/200 [13:13<01:52, 8.02s/pipeline]Optimization
Progress: 94%|██████████| 187/200 [13:17<01:30, 6.96s/pipeline]Optimization Progress: 94%|██████████| 188/200 [13:20<01:05,
5.48s/pipeline]Optimization Progress: 94%|██████████| 189/200 [13:26<01:01, 5.63s/pipeline]Optimization Progress: 95%|██████████|
190/200 [13:28<00:45, 4.57s/pipeline]Optimization Progress: 96%|██████████| 191/200 [13:33<00:42, 4.68s/pipeline]Optimization Prog
ress: 96%|██████████| 192/200 [13:37<00:36, 4.59s/pipeline]Optimization Progress: 96%|██████████| 193/200 [13:39<00:26, 3.77s/pip
eline]Optimization Progress: 97%|██████████| 194/200 [13:43<00:23, 3.84s/pipeline]Optimization Progress: 98%|██████████| 195/200 [
13:45<00:16, 3.33s/pipeline]Optimization Progress: 98%|██████████| 196/200 [13:47<00:12, 3.07s/pipeline]Optimization Progress: 9
8%|██████████| 197/200 [13:50<00:08, 2.84s/pipeline]Optimization Progress: 99%|██████████| 198/200 [13:54<00:06, 3.20s/pipeline]O
ptimization Progress: 100%|██████████| 199/200 [13:59<00:03, 3.68s/pipeline]Optimization Progress: 100%|██████████| 200/200 [14:04<
0:00, 4.28s/pipeline]
Generation 1 - Current best internal CV score: 0.8388888888888889
Optimization Progress: 100%|██████████| 200/200 [14:17<00:00, 4.28s/pipeline]

Best pipeline: GaussianNB(ExtraTreesClassifier(input_matrix, bootstrap=False, criterion=gini, max_features=0.5, min_samples_leaf=1, m
in_samples_split=3, n_estimators=100))
>>>
```



IV. RESULT AND DISCUSSION

For validating the developed model, the dataset has been divided into 70% training and 30% testing subsets. While the Random forest classification, K-nearest neighbor and Tpot models were used to predict the Water Quality Prediction, the RFC, KNN and Tpot were utilized for the water quality classification prediction.

V. CONCLUSION

In this paper, the performance of Random forest classification, K-nearest neighbor in supervised learning and Tpot in AutoML were evaluated to predict the water quality. To this end most dataset related well-known components, such as pH, Conductivity, hardness etc., In the Comparison of Python and AIFramework, Python is predict more accuracy level than the AI framework.

REFERENCES

- [1]. Efficient Prediction of Water Quality Index (WQI) Using Machine Learning Algorithms. Author: Md. Mehedi Hassan, Laboni Akter³, Md. Mushfiqur Rahman⁴. Year: 10 December 2021
- [2]. Efficient Water Quality Prediction Using Supervised Machine Learning Author: Umair Ahmed, Rafia Mumtaz, Hirra Anwar, Asad A. Shah. Year: 24 October 2019
- [3]. Predictive Analysis of Water Quality Parameters using Deep Learning Author: Acharana Solanki, Himanshu Agrawal, Kanchan Kare. Year: 9 September 2015
- [4]. Analysis and prediction of water quality using deep learning and auto deep learning techniques. Authors: D. Venkata Vara Prasad, Lokeswari P. Senthil Kumar. Year: 20 January 2021
- [5]. Automating water quality analysis using ML and auto ML technique Authors: G. Prasanna Medha, S. Harshana, S. Jahnvi Srividya. Year: 20 July 2021
- [6]. Water Quality Prediction Using Artificial Intelligence Algorithms Authors: Fahd Abd Algalil. Year: 30 Dec 2020
- [7]. Predicting and Analyzing Water Quality using Machine Learning: A Comprehensive Model Author: Yafra Khan, Chai Soo See. Year: 21 Nov 2019
- [8]. Water Quality Prediction Method Based on LSTM Neural Network Author: Yuanyuan Wang, Jian Zhou, Kejia Chen. Year: 12 Dec 2017
- [9]. Water quality prediction and classification based on principal component regression and gradient boosting classifier approach Author: Md. Saikat Islam Khan, Nazrul Islam, Sifatul Islam. Year: 3 June 2021

- [10]. River water quality index prediction and uncertainty analysis: A comparative study of machine learning models
Author: Seyed Babak Haji Seyed Asadollah, Ahmad Sharafati, Davide Motta Year: 24 January 2021
- [11]. Ahmed et al., 2019 Umair Ahmed, et al. Efficient water quality prediction using supervised Machine Learning Water, 11 (11) (2019), p. 2210
- [12]. Seyed Babak Asadollah, Haji Seyed, et al. River water quality index prediction and uncertainty analysis: A comparative study of machine learning models Journal of Environmental Chemical Engineering, 9 (1) (2021), Article
- [13]. Ali Najah Ahmed, et al. Machine learning methods for better water quality prediction
- [14]. Shweta Tyagi, et al. Water quality assessment in terms of water quality index American Journal of Water Resources, 1 (3) (2013), pp. 34-38
- [15]. Xiaoping Wang, Fei Zhang, Jianli Ding Evaluation of water quality based on a machine learning algorithm and water quality index for the Ebinur Lake Watershed, China Scientific Reports, 7 (1) (2017), pp. 1-18

BIBLIOGRAPHY



Kalaivazhi Vijayaragavan obtained her Bachelor's degree in Electronics and Communication Engineering from REC – Trichy (Presently NIT – Trichy) in 1996 and her Master's Degree in Computer Science and Engineering from SASTRA University, Thanjavur in 2003. She is currently working as an Associate Professor in the Department of Information Technology at Anjalai Ammal Mahalingam Engineering College, Kovilvenni. She obtained her Doctoral Degree from Anna University, Chennai. Her areas of interest include Machine Learning, Data Analytics, IoT Mobile Communication and Computer Networks.



N. Praveen, Pursuing B.Tech – Information Technology (IT) Final year in Anjalai Ammal Mahalingam Engineering College, THIRUVAR



M.V. Sudharsan, Pursuing B.Tech – Information Technology (IT) Final year in Anjalai Ammal Mahalingam Engineering College, Thiruvarur



P.S. Vijayan, Pursuing B.Tech – Information Technology (IT) Final year in Anjalai Ammal Mahalingam Engineering College, Thiruvarur