

Chronic Liver Disease Detection System based on ML Algorithms

G.M.Padmaja¹, M.Gayatri², K.Raviteja³, M.Tarun⁴, K.Raviteja⁵
¹Assistant Professor, Raghu Institute of Technology, Visakhapatnam, Andhra Pradesh, India.
^{2,3,4,5}Student, Raghu Institute of Technology, Visakhapatnam, Andhra Pradesh, India.

Submitted: 01-06-2022

Revised: 05-06-2022

Accepted: 08-06-2022

ABSTRACT: Chronic liver disease is one of the leading causes of death worldwide and affects many people throughout the world. This chronic liver disease is caused by a variety of factors. We know many of them, such as undiagnosed hepatitis, obesity, and also alcohol misuse. Causes of abnormal nerve function, coughing or vomiting blood, hepatic encephalopathy, kidney failure, liver failure, jaundice, and many more are just symptoms of what's more to come. Diagnosis of this disease is also quite expensive and at the same time, it is very complicated. Therefore, the goal of this paper was to evaluate the performance of various machine learning algorithms to reduce the high cost of diagnosing this chronic liver disease by predicting the disease. In this paper, we used three machine learning algorithms: i.e. Random Forest, XGBoost, and Extra Tree Classifier. The performance of these classification techniques was evaluated against different measurement techniques such as precision, accuracy, f1-score, and recall. Furthermore, our current study's only main focus was on using clinical data to predict liver disease and exploring different ways to represent these data through analysis.

KEYWORDS: Chronic, Machine Learning, Classification, Extra Tree Classifier, XgBoost, Random Forest.

I. INTRODUCTION

The liver is the largest organ in the body, it is needed to digest food and remove toxins from the body. Alcohol consumption and viruses lead to liver damage and sometimes they may also lead to life-threatening conditions. Many types of liver diseases include hepatitis, cirrhosis, liver tumors, liver cancer, and many more. Among them, liver disease and cirrhosis are the main causes of death. Therefore, this liver disease is one of the major health problems in the world and that is the reason

why it is called a chronic disease. Every year, more than 2 million people around the world die from liver disease. According to the Global Burden of Disease report, which was later published by BMC Medicine, in the year 2010, over one million people died from cirrhosis disease which is one of the symptoms of liver disease, and one million from Liver-Cancer. Machine learning has had a significant impact in the medical field for predicting and giving diagnosing methods for liver disease.

Now, machine learning guarantees improved detection and prediction of diseases that have aroused interest in the biomedical field and also increases the objectivity in decision-making. By using these ML techniques, medical problems can be solved easily and the cost of diagnosis can be reduced. In this paper, the main goal was to predict outcomes more efficiently and reduce diagnostic costs in the medical field. In this paper, only three machine learning techniques were applied: RF, XGB, and XTC. The performance of these techniques has been estimated from various aspects such as accuracy, recall, f1-score, and precision.

II. LITERATURE SURVEY

Norziah et al. [1] predicted a hepatitis prognosis disease using the Support Vector Machine (SVM) and Wrapper Method. First, for the classification process, they've used wrapper methods to remove the noise features and let the SVM carry out feature selection to get better accuracy. Feature selection was implemented to minimize noisy or irrelevant data. They have achieved the target result by the combination of both the Wrapper Method and SVM techniques. Accuracy: 81%

Dayanand et al. [2] have predicted three major liver diseases: Liver cancer, Cirrhosis, and Hepatitis with the help of distinct symptoms. They used their project using the Naïve Bayes algorithm and SVM algorithm. When a comparison is done between these two algorithms, it has been done based on their classification accuracy measure. From the experimental results, they concluded that the NB algorithm was the better algorithm that predicted diseases with maximum classification accuracy than the other algorithm. Accuracy: 55.6%

Rakshit et al. [3] have used a computing technique combined with the intelligent diagnosis to detect chronic liver disease which is based on implementing the classification and the type detection for this project in the most complex way. For the project, they have used the Artificial Neural Network (ANN) classification algorithm. Accuracy: 79%

III. PROBLEM STATEMENT

This paper aims to analyze the dataset of liver disease patients, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database and also this dataset, male patients were more when compared to the females. In many machine learning papers, it contains only one type of algorithm. But here, we have used only three types of algorithms and compare all of them for the best result.

IV. METHODOLOGY

Our main goal of this paper is to find a suitable Machine Learning technique that can detect Chronic Liver Disease with high accuracy.

Although we have taken only three Machine learning techniques, (i.e. Random Forest, XgBoost, and eXtra Tree Classifier) we have also used a few more popular ML algorithms in our sample tests, in which these 3 ML algorithms yield high results for the dataset that we have obtained.

From the dataset, we have removed inconsistent data, Encoded the categorical data, and split the data for further analysis. For the execution part, we have used python programming language and the editor used was Google Colab as it is useful for collaborative work and helped us to execute the code easily without any directly installing libraries.

4.1 System Architecture

A system architecture is a conceptual model that defines the structure, behavior, and view of a system. Figure 4.1 describes the behavior and concept that we are using in this paper.

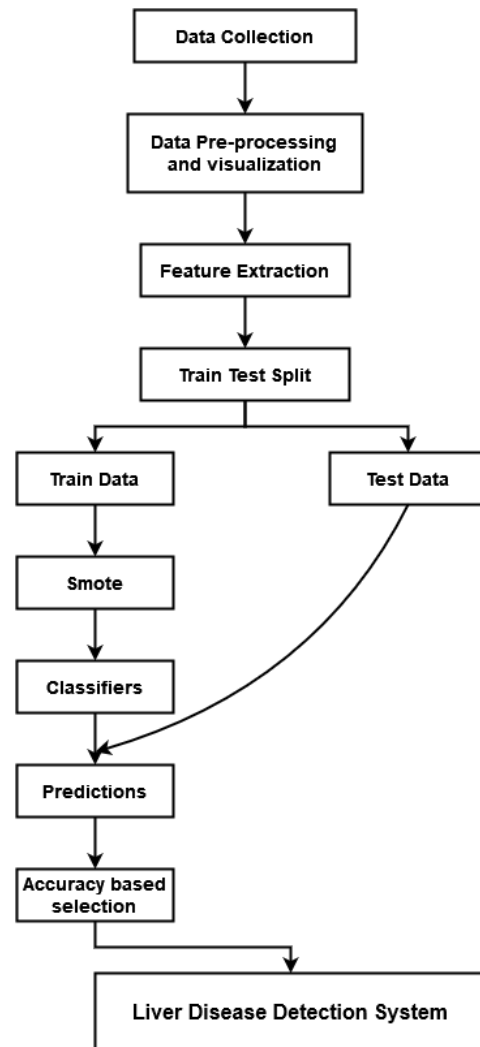


Figure 4.1: System Architecture

4.2 Data Collection

In this paper, we acquired a dataset that is made available by the National Institute of gastroenterology and

Liver Diseases (NIGL D) which has the survey records of Liver disease patients. These survey records were processed by the National Health Portal and were uploaded to various public websites.

This dataset that we have used consists of 583 liver patients' data whereas 75.64% were male patients and 24.36% were female patients also any patient whose age exceeded 89 is listed as being of age "90". This dataset that we have used consists of over 11 particular parameters whereas we chose only ten parameters for our further analysis and one parameter as a target class. The "Dataset" column is a class label used to divide groups into the liver patient (liver disease) or not (no disease). The parameters in the dataset were listed in Table 4.2.

Table 4.2: All attributes in the given Dataset

```

RangeIndex: 583 entries, 0 to 582
Data columns (total 11 columns):
#   Column                               Non-Null Count  Dtype
---  -
0   Age                                   583 non-null    int64
1   Gender                               583 non-null    object
2   Total_Bilirubin                       583 non-null    float64
3   Direct_Bilirubin                      583 non-null    float64
4   Alkaline_Phosphotase                  583 non-null    int64
5   Alamine_Aminotransferase              583 non-null    int64
6   Aspartate_Aminotransferase            583 non-null    int64
7   Total_Protiens                        583 non-null    float64
8   Albumin                               583 non-null    float64
9   Albumin_and_Globulin_Ratio            579 non-null    float64
10  Dataset                               583 non-null    int64
dtypes: float64(5), int64(5), object(1)
memory usage: 50.2+ KB

```

4.3 Datasplitting

After handling the noisy or inconsistent data in the dataset, we split our data into a training set and a testing set. Most of the time we split our data into either 70:30 or 80:20 ratios. For this project, we decided to divide our dataset into an 80:20 ratio and which means eighty percent of the dataset is training data and the remaining twenty percent of our data is testing data.

V. ALGORITHMS USED

5.1 RANDOM FOREST (RF)

Random forest or the random decision forest is one of the ensemble learning techniques of Machine learning for classification, regression, and various assignments that works by spanning a large number of decision trees at the time of training and generating class as a method of classes or a predictive average of each of individual trees. Random decision forests adapt decision trees' propensity to adapt to their given training set. From these combined trees, there is an immediate connection between the combined trees and the results they can achieve.

5.2 eXtreme Gradient Boosting (XGBoost)

XGBoost or eXtreme Gradient Boosting algorithm is one of the popular and efficient open-source algorithms which helps in the implementation of the enhanced-gradient tree algorithm from Classification techniques of Machine-learning. The Enhanced-Gradient tree is a

supervised machine learning algorithm that can attempt to accurately predict a target variable by combining some set of the estimates from a set of a simpler and a weaker model.

5.3 Extra Tree Classifier (XTC)

Extra tree classifier is one of the ensemble learning techniques which was helping in the aggregation of various results of huge uncorrelated decision trees which were collected in a "forest" to generate its classification results. Conceptually, it is very similar to a random forest classifier and differs only in the way decision trees are constructed in the forest. Each decision tree in the additional tree forest is built from the original training samples given to it. Then, at each test node, every tree receives a random sample of features from the feature set from which each decision tree should choose the best feature of all to fit the data according to some mathematical criterion.

VI. RESULTS & ANALYSIS

In this paper, we considered different analyses to examine the three machine learning classifiers for the classification of chronic Liver-Disease datasets. In terms of metrics, i.e. accuracy, recall, f1-score, and precision, the Extra Tree Classifier algorithm (XTC) that we have used has achieved the highest accuracy of over 83.8%, precision: 0.86, recall: 0.83, and f1-score: 0.84 whereas, XGBoost has achieved the least performance among all three ML models used which has an accuracy of 77.8% for the dataset that we have given.

According to the comparison of these measurement criteria, the XTC (Extra Tree Classifier) classification technique in Machine-Learning is much more effective than the other classifiers that we used for predicting this chronic Liver disease for our given dataset. The performance comparison of the three supervised machine-learning techniques that we have used will be shown in the result table (Table 6.1)

Table 6.1: Results of the Algorithms

	Accuracy	Precision	Recall	F1- Score
Random Forest	81.4%	0.84	0.81	0.82
eXtreme Gradient Boosting	77.8%	0.81	0.76	0.78
Extra Tree Classifier	83.8%	0.86	0.83	0.84

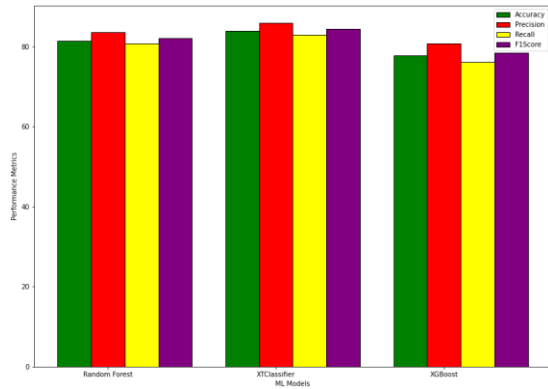


Figure 6.1 Graphical Representation

VII. CONCLUSION

The main objective of this paper is for us to make an effective and an efficient detection system for liver disease patients by using three distinctive and supervised ML techniques where we researched all classifier's execution of all patients' information along with parameters and the Extra Tree classifier gives the most elevated order exactness of 83.8 % measure to predict liver disease and XGBoost only gives the least precision 77.8 %. This project has the option to predict liver infection before advising the well-being condition.

This paper is very surprisingly gainful in low-salary nations where if any absence of medicinal foundations and just particular specialists occurs. We only explored some popular supervised machine learning algorithms, more algorithms can be picked to assemble an increasingly precise model of liver disease prediction and performance can be progressively improved. Additionally, this works likewise ready to assume a significant role in health care research and just as restorative focuses to anticipate liver infection.

REFERENCES

- [1]. A Prediction of Hepatitis Prognosis Using Support Vector Machine and Wrapper Method written by Roslina and A. Norziah, IEEE 2209-22.
- [2]. A paper on Liver Disease Prediction using SVM and NB Algorithms written by S. Dhayanand and S. Vijayarani in the International Journal of Science, Engineering, and Technology Research Vol 4, Issue 4, April 2015.

- [3]. A paper on liver disease detection systems using Machine Learning Techniques written by D. B. Rakshith, S. P. Gururaj, Ashwani Kumar, and Mrigank Srivastava International Journal of Engineering Research & Technology, Vol. 10 Issue 06, June-2021.
- [4]. The dataset that we acquired from Kaggle: kaggle.com/datasets/uciml/indian-liver-patient-records.
- [5]. Types, Symptoms, and Causes of Chronic disease were mentioned detailed in Cleveland Clinic my.clevelandclinic.org/health/diseases/17179-liver-disease.