# A Comparative Study on Various Machine Learning Approaches for the Detection of Alzheimer Disease

## Anuradha Vashishtha[1], Anuja Kumar Acharya[2], and Sujata Swain[3]

*Abstract:* Alzheimer's disease is the most common cause of Dementia. It accelerates the degeneration of brain cells and the progression of memory loss. Identifying and predicting Alzheimer's disease in its early stages is hard. A machine learning system that can predict the disease can solve this problem. High-dimensional data analysis is a massive problem for engineers and researchers in Machine learning. A simple and effective way to solve this problem is to use feature selection to eliminate redundant and useless data. This study aims to determine how accurate different machine learning methods are at diagnosing Alzheimer's disease with help of the feature selection method. Therefore, we utilize the Open Access Series of Imaging Studies dataset in several Machine learning models to accurately detect and predict Alzheimer's disease. We use the wrapper feature selection method in the proposed work to choose the minimal attribute set from the textual records. These approaches include Random Forest and Support vector machine as well as Decision Tree, XGBoost and Voting. We also utilize AdaBoost and Gradient Boosting on the selected data to classify Alzheimer's disease from the Alzheimer's disease vs. Normal disease dataset. As part of the analysis of the results, we found that the machine learning with feature selection methods provides better result as compared to machine learning without feature selection methods.

*Keywords:* Alzheimer's disease (AD); Feature Selection (FS); MRI; Machine Learning (ML); Classification.

## 1. Introduction

Alzheimer's disease is the most frequent cause of Dementia in the world today. An irreversible and inherited brain disease, Alzheimer's syndrome impairs a person's ability to do daily chores, recall information, and engage in higher-level thinking skills. Alzheimer's disease happens when many neurons lose their synaptic connections. Alzheimer's disease is rare in adults between 30 and 60. Alzheimer's disease is characterized by changes in sleep patterns, depression, anxiety, and trouble with simple tasks like reading and writing [1]. It is thought that by 2050, Dementia, which affects about 40–50 million people now, will affect more than 131.5 million people worldwide [3]. About 70% of people in low-income countries have Dementia.

Alzheimer's disease usually gets progressively over time in three stages: early, middle, and late. In medical terms, these stages are sometimes called mild, moderate, and severe. When Alzheimer's is in its early stages, a person may still be able to live on their own. They can still drive, work, and do things with friends and family [4]. As a result, the individual may feel they are suffering from memory issues, such as forgetting what they have just read or having difficulty performing duties in social or work contexts. It is during the middle stage of Alzheimer's disease that people typically experience the most debilitating symptoms. People

around you can witness your symptoms, such as forgetting recent events or your past. Since they could not remember their address, phone number, or where they went to high school or college, Changes in personality and behavior, such as being suspicious, having delusions, or doing things like wringing their hands or shredding paper repeatedly. In the last stage of the disease, the person may need help with daily tasks and personal care around the clock, have trouble communicating, and see changes in their physical abilities, like walking, sitting, and eventually swallowing. High-dimensional data analysis is a big problem for engineers and researchers in ML. Applying FS to eliminate redundant and useless data is a simple and effective way to solve this problem. Getting rid of the useless data improves learning accuracy, shortens the time it takes to do calculations, and makes it easier to understand the learning model or data. Most of the time, not all the variables in a dataset are helpful when building an ML model in the real world. When you add redundant variables to a model, it becomes less good at generalization and may also be less accurate.

Further, if you add more variables to a model, you end up with a more complicated model. To overcome this issue, we use the feature selection method to select the most relevant elements from a text-based database. In ML, the goal of FS is to find the best set of features that can be used to make applicable models of what is being studied [5]. In ML, the ways to choose which features to use can be put into the following types:

### 1.1. Supervised methods

These methods are used for labeled data and to group the essential features to make supervised models like classification and regression work better.

[1] *Kalinga Institute of Industrial Technology (KIIT), deemed to be University, Bhubaneswar – 751024, ODISSA*
*ORCID ID :  0000-0003-1270-4139*

[2] *Kalinga Institute of Industrial Technology (KIIT), deemed to be University, Bhubaneswar – 751024, ODISSA*
*ORCID ID :  0000-0002-1673-3378*

[3] *Kalinga Institute of Industrial Technology (KIIT), deemed to be University, Bhubaneswar – 751024, ODISSA*
*ORCID ID :  0000-0001-7089-1863*
*\* Corresponding Author Email: anuradhavashishtha30@gmail.com*

## 1.2. Unsupervised methods

When there are no labels, these methods are used. From the point of view of taxonomy, these methods can be put into:

### 1.2.1. Filter methods:

The essential qualities of features, which can be quantified with univariate statistics, are collected instead of cross-validation results. When working with high-dimensional data, using filter methods is cheaper in computing.

### 1.2.2. Wrapper methods:

Wrappers need a way to search for all possible subsets of features, test a classifier with that subset of parts, and figure out how good that Classifier is by learning. The process of choosing which features to use is based on an ML algorithm that is tried to fit a given dataset. Most of the time, wrapper methods are better at making predictions than filter methods.

### 1.2.3. Embedded methods:

These methods combine the best parts of both filter and wrapper methods. They consider how features interact with each other and keep a reasonable cost of computation. It might be hard to pick out a small set of features.

Machine learning (ML) is the process of discovering procedures, building models, and then applying what it has learned to produce results on its own [7]. A model is a machine learning system that has been taught to recognize certain kinds of patterns by using an algorithm. That means it works with the data and finds the structures that are hidden in the data. The known answers and extracted features of a dataset are used to come up with a formula that uses the input and output functions to predict the response. So, the model's algorithm takes a training set of data to construct a strategy to predict the output and then saves that approach for future usage. It has an algorithm that automatically pulls out the data and uses data mining to build models from the data to predict what will happen in the future. CAD stands for "computer-aided detection." It is a type of computer system that helps doctors make more accurate diagnoses [8]. When a CAD system is combined with machine learning (ML), it gives healthcare systems a good way to turn those data into useful information. This makes it easier to find AD early and accurately.

In our work, we look for people who have Alzheimer's, and we try to find people who might have Alzheimer's at an early stage. Both OASIS and Kaggle have datasets for Alzheimer's disease, which are used to select the minimal attribute for the prediction of AD. It is also used to train all the patient data using different ML algorithms, such as Random Forest classifier (RF), SVM, Decision tree classifier (DT), XGBoost, voting classifier, and other Classifier to be able to tell the affected people apart quickly and efficiently.

Here are the different parts of our study: In Section 1, we talk about what is Alzheimer's disease. In section 2, what has been done recently to find it using Machine Learning models using feature section and without feature selection. In Section 3, we talk about data description, the feature selection method, and different Machine Learning classifier models. Section 4, a comparison of various machine learning models in the dataset. Section 5, the conclusion of the work and talks about what needs to be done in the future.

## 2. Related Work

This section provides a detailed literature review of the various approaches for AD and NC classification. The existing literature on AD and NC is classified based on various machine learning approaches. Bansal, Deepika, et al. [9] using various machine learning algorithms for Alzheimer's disease classification. They are using the oasis dataset, and the dataset includes 416 participants in cross-sectional data and 373 entries in longitudinal data. In order to remove superfluous attributes, CFSSubsetEval is employed. The individual predictive ability of each feature, as well as the degree of redundancy among them, is used to evaluate a feature subset. J48, Random Forest, Nave Bayes, and Multilayer Perceptron are the four classifiers employed. Shahbaz, Muhammad, et al. [10] used machine learning and data mining to improve the quality and efficiency of medical care. Six machine learning and data mining algorithms, such as k-nearest neighbors (k-NN), rule induction, Naive Bayes, generalized linear model (GLM), and deep learning algorithm, were used on the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset to classify the five stages of AD.

Leong, Lee Kuok, and Azian Azamimi Abdullah [11] work is based on comparing and evaluating different machine learning algorithms with pre-processing and Boruta algorithm feature selection in the diagnosis and prediction of AD. In the suggested method, the missing data were taken out by taking out columns from the SES and MMSE columns. Pre-processing and experimental results with Boruta characteristics selection is a good way to find Alzheimer's disease early. With an accuracy of 94.39 percent, RF GSCV with the Boruta algorithm is the best Classifier out of all the ones that are available.

Alickovic, Emina, and Abdulhamit Subasi [13] different types of machine learning have been widely utilized to diagnose Alzheimer's disease. Alickovic and Subasi's study shows that supervised machine learning models can be used to predict Alzheimer's disease. The most important parts of this research were the support vector machine, naive Bayes, k-nearest neighbors, random forest, artificial neural network, and logistic regression. They did their research with the help of the ADNI data repository (ADNI, 2017). The most accurate classifiers are the random forest classifier, which is 85.77% accurate, and the k-nearest neighbors' Classifier, which is 84.27 % accurate.

Deep convolutional autoencoders are used by Martinez-Murcia et al. [14] to investigate the data analysis of AD. MRI images can be broken down using data to find a person's cognitive symptoms and the neurodegenerative process that is causing them. A regression and classification analysis and an estimate of the effect of each coordinate on the brain manifold are used to find out how the extracted features are spread out in different ways. Scores on the MMSE or ADAS11, along with signs from imaging, can be used to predict with over 80% accuracy a diagnosis of Alzheimer's disease.

Alroobaea, Roobaea, et al. [15] proposed using ML to detect Alzheimer's. ADNI and OASIS are used to evaluate categorization models. The experimental results showed that logistic regression and support vector machine offered the best accuracy using the ADNI dataset, while logistic regression and random forest supplied the best accuracy using the OASIS dataset. Sinan Faouri et al. [16] research examines machine learning algorithms for predicting Dementia in elderly adults. The dementia dataset uses seven machine learning techniques and two feature reduction methods with distinct selection thresholds. This experiment's

results are: Information gain seems more efficient and stable than PCA for ranking and selecting features. Machine learning algorithms are not stable across all accuracy measures. The most stable ways for machines to learn are SVM and NB. In PCA, it is better to use a lot of dimensions than a few. The accuracy of group learning methods like Ada and Random Forest was not stable. By raising the information gain threshold, adding more features to solo models made them better, though some features may not be needed.

Kavitha, C., et al. [17] In this research, they identify Alzheimer's patients and look for early-stage cases. Alzheimer's datasets are available on OASIS and Kaggle, and they are used to train all patient data using machine learning algorithms like SVM, Random Forest classifier, Decision tree classifier, XGBoost, and Voting classifier to effectively tell affected people apart quickly and efficiently. Khan, Afreen, and Swaleha Zubair [18] came up with a five-step machine learning pipeline, with each step being further broken down. The Open Access Series of Imaging Studies (OASIS) database of MRI brain images was used to do the analysis for the study. One hundred fifty people had a total of 343 MRI appointments. Scores from the MMSE, CDR, and ASF (Atlas Scaling Factor) were used in the analysis. The proposed ML pipeline has a classifier system, a way to change data, and a way to choose features. All these parts are wrapped in a framework for experimenting and analyzing data.

It is observed from the literature that they used the various ML method for AD classification, have a drawback: the lack of willingness to choose features and search efficiently to find the sub-optimal features for classifying AD and NC. So, it is important to use a feature-selection algorithm that can figure out what the most important things are that help classify AD and NC.

## 3. Material and methods Units

### 3.1. Dataset

The workflow is built with MRI data from the Open Access Series of Imaging Studies (OASIS) database. We focused our research on getting MRI data over time from both older people with Dementia and older people without Dementia. Researchers can use MRI datasets that are stored in Oasis, which is an open-source database.

### 3.1.1. Details of the acquitted MR images

All the structural MRI images were taken with the 1.5 Tesla Vision Scanner. The high-resolution MP-RAGE (Magnetization Prepared Rapid Acquired Gradient Echo) sequence was used to look at 343 MRI scans of 150 different people [19]. Table 1 has information about how to get an MRI.

**Table 1.** Details of the MRI acquisition

| MR characteristics | Values |
|---|---|
| TR (repetition time) | 9.7 msec |
| TE (echo time) | 4.0 msec |
| Flip angle | $10^{\circ}$ |
| T1 | 20 msec |
| TD | 200 msec |
| Orientation | Sagittal |
| Thickness | 1.25 mm |
| Gap | 0 mm |
| Slice number | 128 |
| Resolution | $256 \times 256 \ (1 \times 1 \ mm)$ |

### 3.1.2. Subjects

A total of 343 sessions were conducted on 150 people, ranging in age from 60 to 96. A longitudinal study of the population is included. These individuals' demographics may be seen in Table 2.

**Table 2.** Details of the MRI acquisition

| No of subjects | 78 Demented (AD) | 72 Non-Demented (ND) |
|---|---|---|
| Male | 40 AD | 22 ND |
| Female | 38 AD | 50ND |
| Age: Range (in years) Means ± SD Median | 60-96 77.01 ± 7.64 77.0 | - |

### 3.1.3. Dataset description

The OASIS has just released OASIS-3, the latest in a series of releases aimed at making neuroimaging datasets openly available to researchers. It is our goal that this multi-modal dataset will aid future research in both basic and clinical neuroscience. Researchers have used previously available OASIS data [20] (from Marcus et al., 2007) and OASIS-longitudinal data (from Marcus et al., 2010) to test hypotheses, create atlases of the brain, and create segmentation algorithms. To better understand aging and Alzheimer's disease, researchers are using the OASIS-3 dataset. It is hoped that the OASIS datasets maintained at central.xnat.org would offer the community open access to a sizable library of neuroimaging and processed imaging data spanning a broad demographic and cognitive and genetic range. All the data may be accessed at www.oasis-brains.org.

The dataset has 373 observations and 15 features. These 15 features names are Subject ID, MRI ID, Group, Visit, MR delay, M/F, Hand, Age, EDUC, SES, MMSE, CDR, eTIV, nWBV, ASF. In the dataset, the goal variable 'Group' is a binary classifier that defines whether a patient has Dementia. In this work, we utilized scores to compare healthy and diseased brains. The Alzheimer's disease in the people who took part in the study ranged from very mild to moderate. All the subjects were right-handed, and they

were both men and women. All the patients went through the same analysis process. All the people in the control group took the Mini-Mental State Examination (MMSE) and other tests as part of a neuropsychological evaluation. All the control subjects were checked out with an MRI using the same machine and method as the AD patients. The dataset includes 373 MRI sessions from non-demented, demented, and converted patients. On their initial appointment, i.e., Visit = 1, some individuals were non-demented and eventually diagnosed with Dementia. Converted patients were thus categorized. This study solely considers patients with Visit = 1. Table 3 indicates the number of patients in each group, totaling 150 participants investigated.

## 3.2. Proposed model

In our proposed model we consist of six steps as shown in Fig 1. In the first step i.e., Data pre-processing step, collecting the data, and handling missing values. In the second step, we are using a feature selection method to select the minimal features from the OASIS MRI dataset. In the third step, data segregation entails dividing the dataset into train and test sets. After that, the data are utilized for the creation of the model in the fourth step. The fourth step consisted of two sub-steps: (i) model training and (ii) model evaluation. It forecasts the model on the test set, categorizing people as having Dementia or Non dementia. The final step, performance evaluation is illustrated the model's performance graphically to provide insights.

### 3.2.1. Data pre-processing

Data pre-processing is an important part of building any machine learning model because it makes sure that the training data is good. In this work, some important steps were taken to prepare the data, such as dealing with missing values, finding outliers, and normalizing the data.

### 3.2.2. Feature selection

In machine learning, it is very important to choose the right features. In this work, feature selection is used on a huge number of samples of clinical data about Alzheimer's disease. There are three ways to choose features: the filter method, the wrapper method, and the embedded method. Wrapper method can be used either using backward feature selection or forward feature selection or exhaustive feature selection. As discussed in the literature review, exhaustive feature selection performs better than other feature selection methods. Keeping this, we have used Exhaustive Feature selection method in the current proposed work.

### 3.2.2.1. Exhaustive Feature selection (EFS)
EFS compares the performance of an ML method to the performance of all possible combinations of the features in the dataset. It is decided to include the feature subset that offers the best level of performance. To find the ideal subset of characteristics, EFS examines all potential feature combinations.

---

**Algorithms 1 Exhaustive feature selection**

1: /* X is a group of features */

2: /* n is the number of features */

3: /* Size of a subset of the target features d */

4: / * Set of all the possible subsets of X's features, F, where each subset has size d */

5: $Y_{opt} \leftarrow \emptyset$

6: $G_{opt} \leftarrow -\infty$

7:

8:

9:

6:

7:

8:

9: end for

---

**Table 3.** Number of patients for each group on the First Visit

| Group (target variable) | Number of patients |
|---|---|
| Non-Demented | 72 |
| Demented | 64 |
| Converted | 14 |

---

### 3.2.3. Machine learning Classifier Models

After feature selection, the data is sent to the various Classifier for the detection of the AD class. Here we have used RF, SVM, Decision tree, XGboost, Voting, extra tree, and others.

### 3.2.3.1. Random Forest (RF)
It is better to use an RF model than a decision tree since it avoids the problem of "overfitting." They all have slightly different decision trees, but they all form RF models. The ensemble uses a majority voting technique to create predictions based on each individual DT model (bagging). So, the amount of overfitting is cut down while each tree's ability to predict stays the same.

### 3.2.3.2. Support Vector Machine (SVM)
In this approach, the class of data points in multidimensional space are identified using appropriate hyperplanes as the determining factor. By using SVM [22], our objective is to locate a hyperplane that distinguishes between two groups of variables that each occupy adjacent clusters of vectors, one on one side and the other on the other side of the hyperplane. The vectors that are closest to the hyperplane are support vectors. In SVM, both training and test data are utilized. Data for training is segmented into the goal values and attributes. The support vector machine (SVM) generates a model for forecasting target values based on test data.

### 3.2.3.3. Decision Tree (DT)
An overview of a DT uses the feature cut-off values to divide the data repeatedly. By dividing a set into subgroups, splitting creates new subsets. Internal nodes are intermediate subsets, and leaf nodes are leaf nodes. When features and the objective have a lot of interaction, a decision tree is a great tool.

### 3.2.3.4. XGBoost
It stands for eXtreme Gradient boosting, which is also written as XGBoost. Here, the goal is to use gradient-boosted decision trees to get the most speed and performance out of them. Because of

this, Gradient boosting machines can be slow to set up and do not scale well. The whole point of XGBoost is to get things done quickly [17].

### 3.2.3.5. Voting

One vote is all that is required to integrate multiple learning algorithms. In order to make use of their unique properties, voting classifiers are not actually classifiers but rather wrappers for several other classifiers that are trained and evaluated simultaneously. We can use data sets to train algorithms and ensembles that can predict the outcome. You can give a prediction vote in two ways:

Hard Voting: Hard Voting is the most basic kind of majority voting. In this case, Nc will win because it has the most votes. We use most of each Classifier to make our predictions.

Using a "soft vote," the probability vectors for each projected class are combined, and the highest value is chosen. Every Classifier is subjected to this procedure (recommended only when the classifiers are well calibrated) [17].
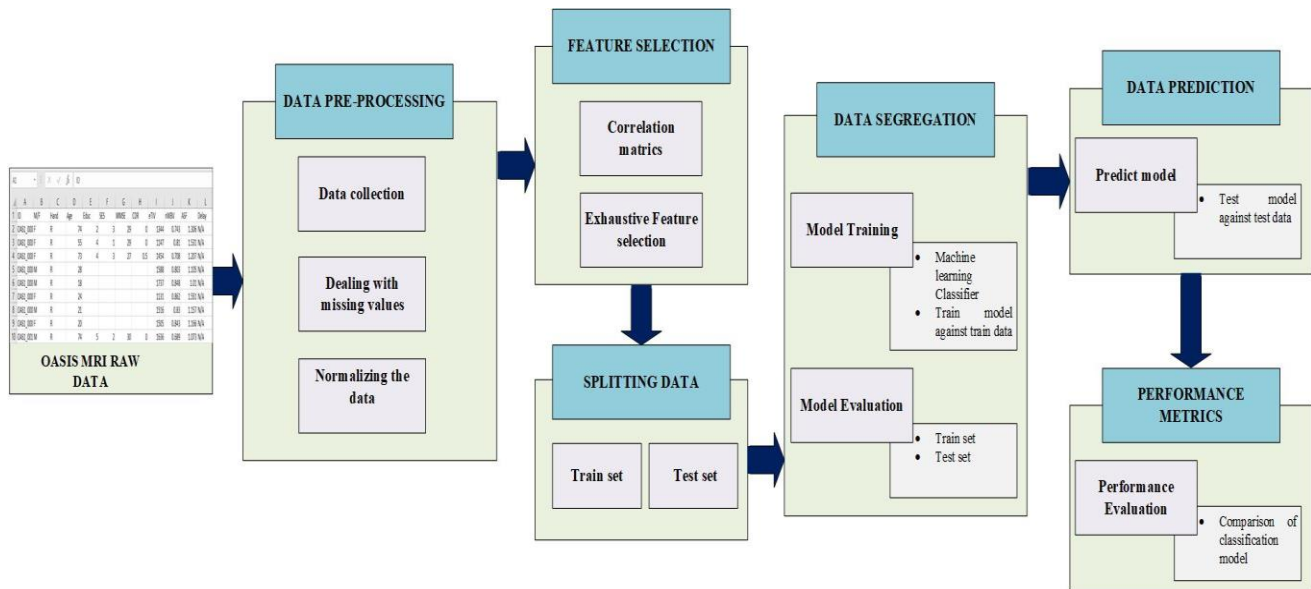


**Fig. 1.** The workflow diagram of the proposed model

### 3.2.3.6. Extra Trees

A massively randomized tree classifier is utilized in ensemble classification methods. Unlike decision tree classifiers, it is created differently. Compared to random forest classifiers, they are also significantly faster. It uses a forest-like structure to collect the outputs of several disconnected decision trees and then combines them to get a categorization result. There are multiple randomized decision trees that are used to fit extra-tress on various dataset attributes [23]. Improve the classification accuracy by limiting over-fitting and using the averaging method. Additionally, the secondary training set is used to construct each of the extra trees in the forest.

### 3.2.3.7. Gradient Boosting

Prediction models that are ensembles of weak prediction models (often decision trees) are generated via Gradient Boosting [24]. It combines numerous poor ML classifiers to build a single powerful one. It takes a step-by-step method to build a model. Furthermore, it offers optimization of loss functions that can be randomized and differed. The AdaBoost classifier is the basis for this Classifier, which is superior to the Adaboosting approach. Both classifiers and weighted inputs are recalculated using an AdaBoost approach combined with an approach known as "weighted minimization." gradient-boosting classifiers are designed to minimize loss, which is the difference between a test set's actual classification and its predicted classification.

### 3.2.3.8. AdaBoost

The focus of this paper is solely on issues of binary categorization. Classifier and estimator for groups. Ensemble classifiers combine multiple machine learning classifiers into one algorithm. Using AdaBoost turns weak ML classifiers into strong Classifiers. First, it is a meta-estimator, which means that it fits the Classifier on the training dataset before moving on to fitting more classifier replicas on top of those on the original train dataset. Instead, the weights that were given to events that were wrongly classified are recalculated in a way that will make future classifiers pay more attention to hard cases [25].

## 4. Experimental Results and Discussion

We are using the Python-based computing package that works in the GPU environment. The experimental work has been put into NVIDIA cards that support up to 8 GB of RAM. In our study, we use a method called feature selection to help reduce the number of attributes and improve accuracy. In order to measure the performance parameters like accuracy, precision, recall, and F1-score are being evaluated for the above-mentioned classifiers.

### 4.1. Performance evaluation

In this section, we evaluate the precision, F1-score, and recall for measuring the performance of various methods.

Accuracy: It is the percentage of correctly categorized results out of all the samples that are used to calculate accuracy. The Equation

(i) is used to figure out the accuracy.

### 4.1.1. Prediction accuracy:

This is the percentage of predicted positive outcomes that are correct. For a decent classifier, the Precision value should be 1. The Equation (ii) is used to calculate the precision.

### 4.1.2. Recall:

Recall (REC) is a rate that is good. If the recall is 1, it means that it is a good sorter. The recall is calculated using the Equation (iii).

### 4.1.3. F1 Score:

It is a measure that considers both Recall and Precision. Only when both measures, like Recall and Precision, are 1 does the F1 score become 1. The precision and recall are shown by the Equation (iv).

$$Accuracy\ (ACC) = \frac{TP + TN}{(TP + TN + FP + FN)} \qquad \text{i}$$

$$Precission\ (PR) = \frac{TP}{TP + FP} \qquad \text{ii}$$

$$Recall\ (REC) = \frac{TP}{(TP + FN)} \qquad \text{iii}$$

$$F1 - Score = \frac{2 * PR * REC}{PR + REC} \qquad \text{iv}$$

In statistical terms, correlation is a way of determining how closely related two variables are to one another [26]. Positivity or negativity are both possible outcomes. A rise in the value of a single feature correlates positively with the goal variable's value, while an increase in a single feature's value correlates negatively with the target variable's value. As a result, the Heatmap was used to generate the ML model's correlation matrix. In order to discover the features that are most closely linked to the dependent variable, a heatmap is a visual representation of data. The removal of associated variables is a key step in model construction. Fig 2 shows the correlation matrix with Heatmap. We are using only 3 relevant features in our proposed work, and that feature is selected by Exhaustive Feature selection. Fig 3 shows the important feature graph).

Most popular metrics include True Positive (TP), False Positive (FP), True Negative (TN), and False Negative (FN) measures. Fig 5, Fig 6, Fig 7, Fig 9 show the confusion matrix for ML models like Random Forest, SVM, Decision tree, XG boost, Soft &Voting Classifier, and others. ROC graphs for RF, SVM, Decision Tree, XG Boost and Soft & Hard Voting classifiers are also shown in Fig 5, Fig 6, Fig 7, Fig 9. An accurate comparison of each machine model is depicted in Fig 10. A detailed examination of all the experiments' outcomes (including the averages of the models' accuracy, precision, recall, and the F measure) was performed, as can be seen in Table 4, Table 4, Table 5, Table 6. The accuracy, precision, recall, and F1 score of each of the Machine Learning models are all compared. Fig 4 shows the ROC curve and a comparison of the AUC for Random Forest, SVM, Decision tree, XG boost, Voting, extra tree, Gradient Boosting, and AdaBoost.
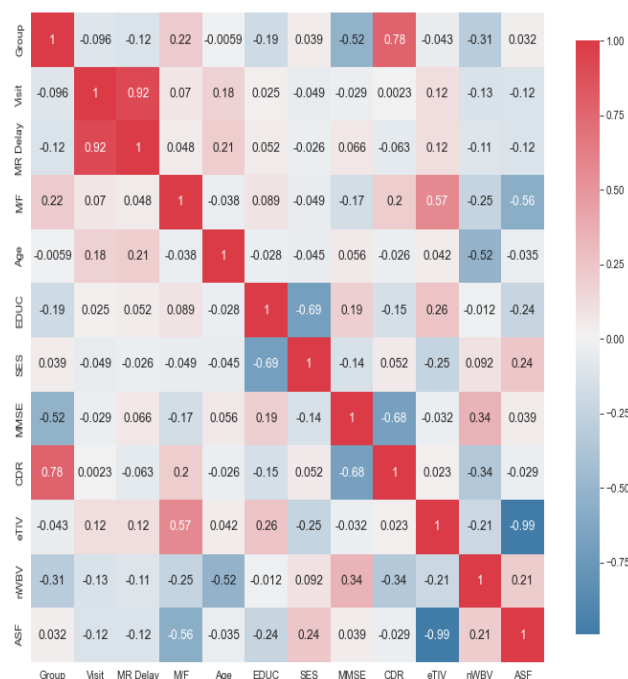


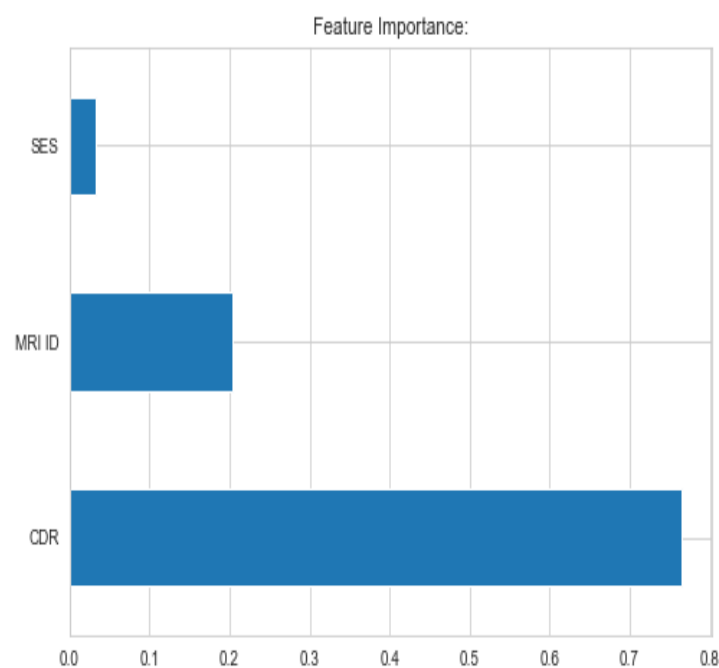**Fig. 2.** Feature selection using correlation matrix with Heatmap.



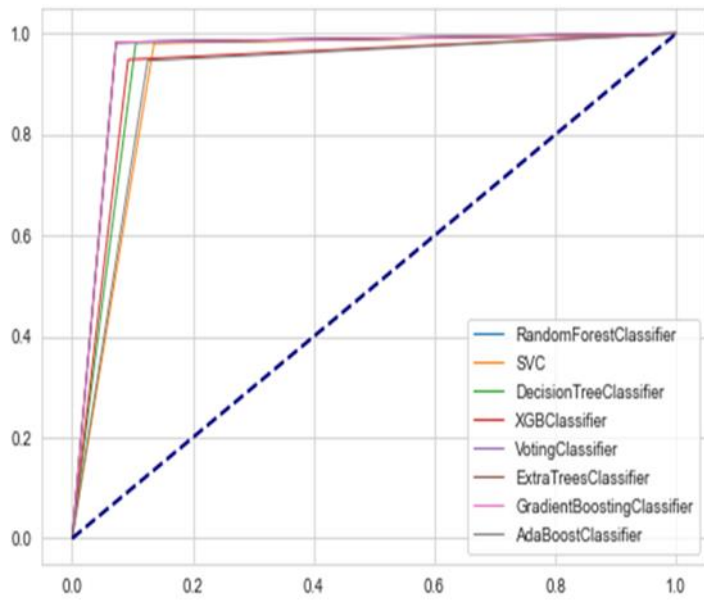**Fig. 2.** Feature selection using feature importance.
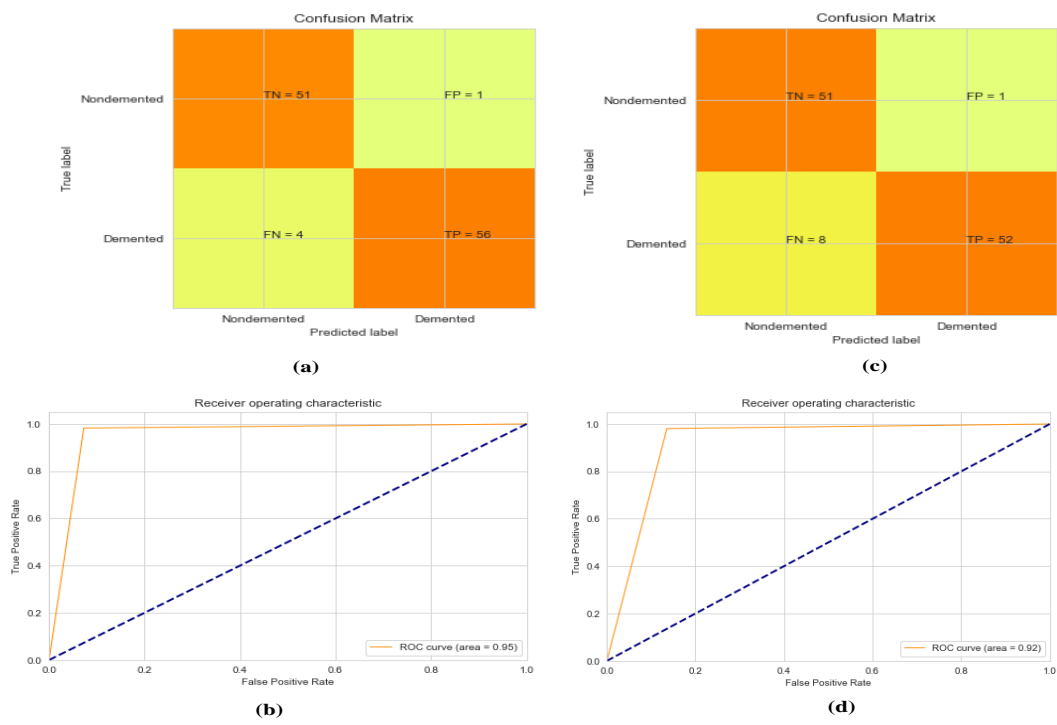
Fig. 4. Plotting of ROC and comparison of AUC



(a)



(c)



(b)



(d)

**Fig. 5.** **(a)** Confusion matrix of random forest. **(b)** ROC curve of random forest **(c)** Confusion matrix of SVM and **(d)** ROC curve of SVM
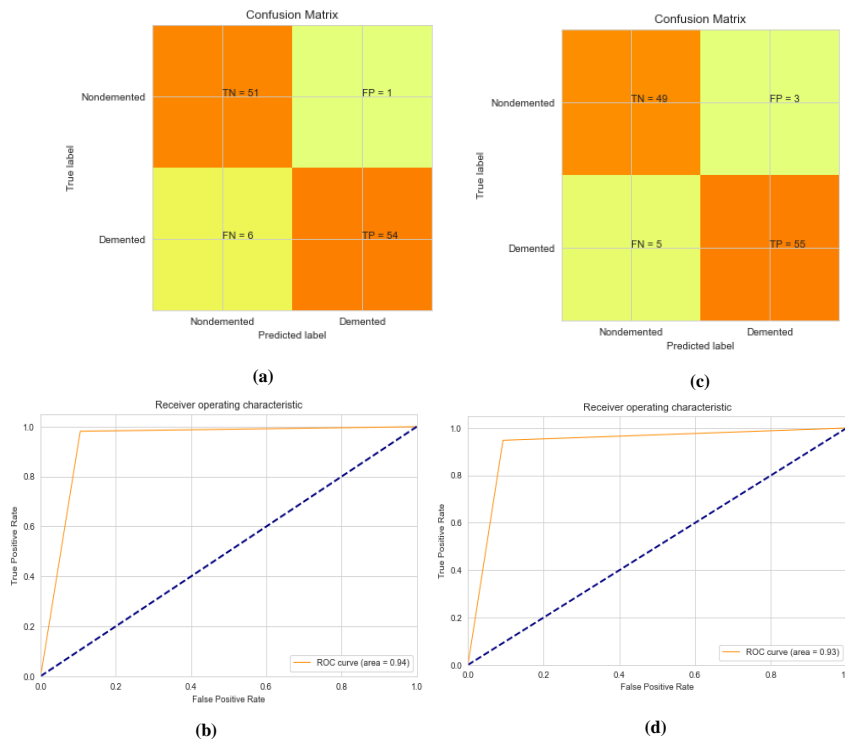
**Fig. 6.** **(a)** Confusion matrix of Decision Tree Classifier. **(b)** ROC curve of Decision Tree Classifier **(c)** Confusion matrix of XGBOOST and **(d)** ROC curve of XGBoost
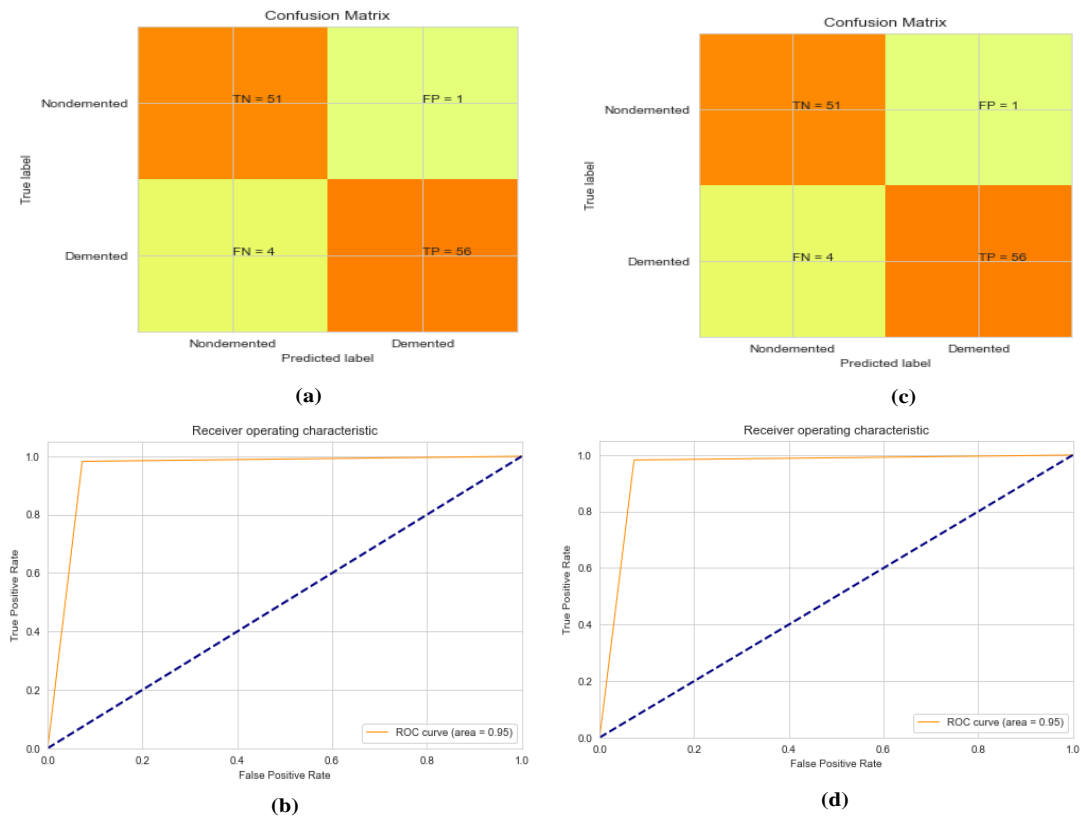


**Fig. 7.** **(a)** Confusion matrix of vote hard Classifier. **(b)** ROC curve of vote hard Classifier **(c)** Confusion matrix of vote soft and (d) ROC curve of vote soft.
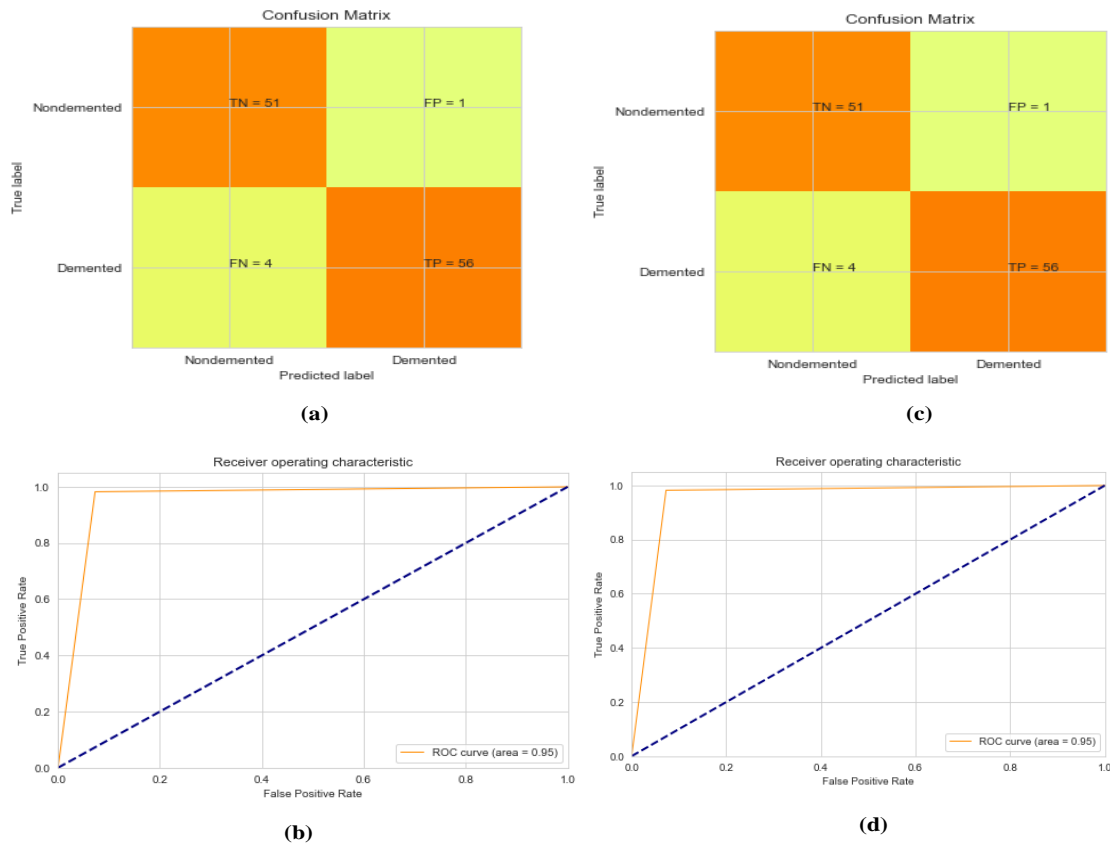
**Fig. 8. (a)** Confusion matrix of Extra Trees Classifier. **(b)** ROC curve of Extra Trees Classifier **(c)** Confusion matrix of Gradient Boosting and **(d)** ROC curve of Gradient Boosting.
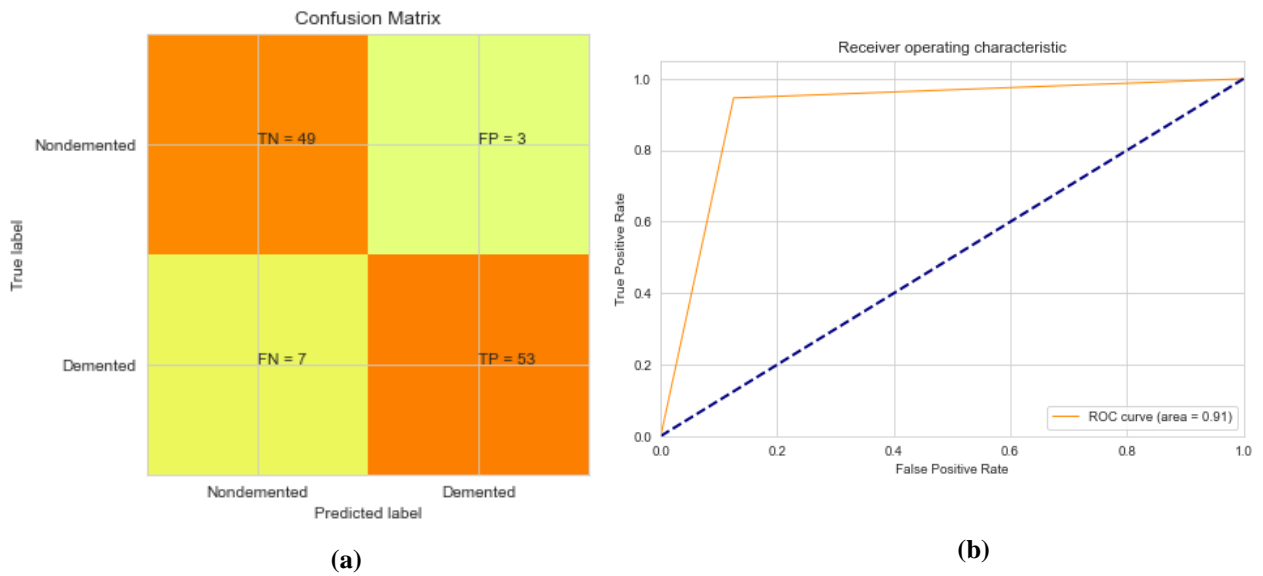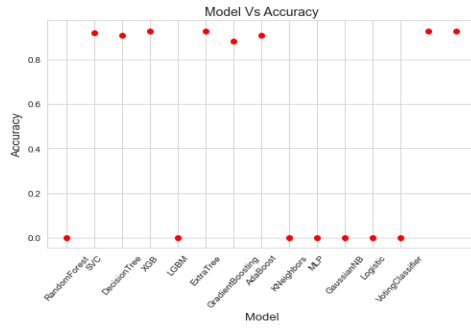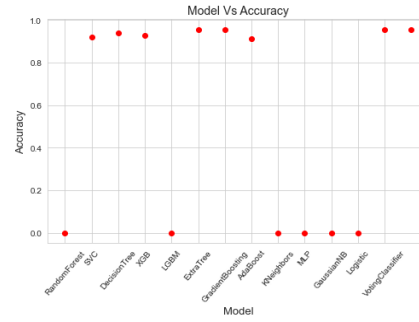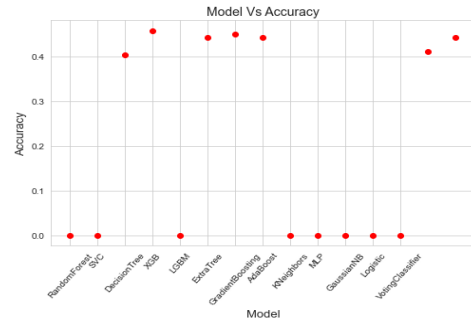


**Fig. 9. (a)** Confusion matrix of AdaBoost Classifier **(b)** ROC curve of AdaBoost Classifier
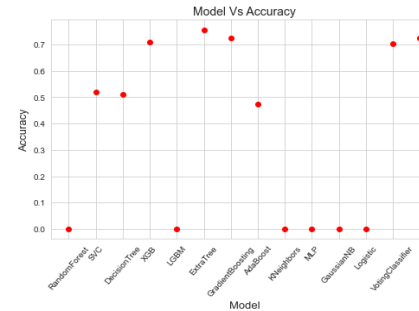
**(a)** Using a longitudinal dataset to compare accuracy without selecting features.



**(b)** Using a longitudinal dataset to compare accuracy with feature selection.



**(c)** Using a cross-sectional dataset to compare accuracy without feature selection.



**(d)** Using cross-sectional to compare the accuracy without selecting features.

**Table 4.** Comparison accuracy table of models without feature selection using longitudinal dataset

| Model | ACC | PR | REC | F1-Score |
|---|---|---|---|---|
| Random Forest | 0.92 | 0.98 | 0.88 | 0.93 |
| SVM | 0.91 | 0.98 | 0.87 | 0.92 |
| Decision Tree | 0.91 | 0.98 | 0.85 | 0.91 |
| XGBoost | 0.92 | 0.96 | 0.90 | 0.93 |
| Voting | 0.92 | 0.96 | 0.90 | 0.93 |
|  | 0.92 | 0.96 | 0.90 | 0.93 |
| Extra Tree | 0.92 | 0.96 | 0.90 | 0.93 |
| Gradient Boosting | 0.88 | 0.91 | 0.87 | 0.89 |
| AdaBoost | 0.91 | 0.93 | 0.90 | 0.92 |

**Table 6.** Comparison accuracy table of models without feature selection using Cross- dataset

| Model | ACC | PR | REC | F1-Score |
|---|---|---|---|---|
| Random Forest | 0.45 | 0.45 | 0.36 | 0.40 |
| SVM | 0.40 | 0.37 | 0.27 | 0.31 |
| Decision Tree | 0.40 | 0.39 | 0.33 | 0.36 |
| XGBoost | 0.45 | 0.46 | 0.42 | 0.44 |
| Voting | 0.43 | 0.43 | 0.36 | 0.39 |
|  | 0.41 | 0.42 | 0.38 | 0.40 |
| Extra Tree | 0.48 | 0.48 | 0.39 | 0.43 |
| Gradient Boosting | 0.44 | 0.45 | 0.44 | 0.44 |
| AdaBoost | 0.44 | 0.41 | 0.23 | 0.29 |

**Table 5.** Comparison accuracy table of models with feature selection using longitudinal dataset

| Model | ACC | PR | REC | F1-Score |
|---|---|---|---|---|
| Random Forest | 0.95 | 0.98 | 0.93 | 0.96 |
| SVM | 0.91 | 0.98 | 0.87 | 0.92 |
| Decision Tree | 0.93 | 0.98 | 0.90 | 0.94 |
| XGBoost | 0.92 | 0.95 | 0.92 | 0.93 |
| Voting | 0.95 | 0.98 | 0.93 | 0.96 |
|  | 0.95 | 0.98 | 0.93 | 0.96 |
| Extra Tree | 0.95 | 0.98 | 0.93 | 0.96 |
| Gradient Boosting | 0.95 | 0.98 | 0.93 | 0.96 |
| AdaBoost | 0.91 | 0.95 | 0.88 | 0.91 |

**Table 7.** Comparison accuracy table of models with feature selection using Cross- dataset

| Model | ACC | PR | REC | F1-Score |
|---|---|---|---|---|
| Random Forest | 0.67 | 0.71 | 0.61 | 0.66 |
| SVM | 0.51 | 0.53 | 0.45 | 0.49 |
| Decision Tree | 0.51 | 0.53 | 0.26 | 0.35 |
| XGBoost | 0.70 | 0.75 | 0.64 | 0.69 |
| Voting | 0.68 | 0.73 | 0.61 | 0.66 |
|  | 0.70 | 0.73 | 0.68 | 0.70 |
| Extra Tree | 0.67 | 0.72 | 0.59 | 0.65 |
| Gradient Boosting | 0.72 | 0.75 | 0.68 | 0.71 |
| AdaBoost | 0.47 | 0.44 | 0.17 | 0.24 |

## 5. Conclusion

Alzheimer's disease is a major health concern, and instead of trying to find a cure, it is more important to lower the risk, start treatment early, and find the symptoms quickly and correctly. As the literature review shows, there have been a lot of attempts to find Alzheimer's Disease using different machine learning algorithms and micro-simulation methods. However, it is still hard to find relevant attributes that can find Alzheimer's very early. In this current study wrapper method found a sub-optimal minimal feature set for classifying the AD and NC patients with better classification performance on the textual data. Moreover, the wrapper method also returns a minimal feature subset for AD and NC classification within less time complexity. The sub-optimal minimal features returned by the wrapper method after executing with the various machine learning classifiers also increased the classification accuracy for AD and NC patients on the dataset. Further research will focus on extracting and analysing novel features that can better aid in the diagnosis and elimination of Alzheimer's Disease from existing feature sets in order to increase detection accuracy.

## 6. References and Footnotes

### Conflicts of interest

There is no conflict of interest, the authors say. The people who paid for the study had nothing to do with how it was set up, how the data were collected, analyzed, or interpreted, how the manuscript was written, or whether or not the results were made public.

### References.

[1] Lodha, Priyanka, Ajay Talele, and Kishori Degaonkar. "Diagnosis of alzheimer's disease using machine learning." 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA). IEEE, 2018.

[2] Yang, J.-P. . "A Novel Storage Virtualization Scheme for Network Storage Systems". International Journal on Recent and Innovation Trends in Computing and Communication, vol. 10, no. 1, Jan. 2022, pp. 08-13, doi:10.17762/ijritcc.v10i1.5514.

[3] Bari Antor, Morshedul, et al. "A comparative analysis of machine learning algorithms to predict alzheimer's disease." Journal of Healthcare Engineering 2021 (2021).

[4] A. Simon, M. Deo, V. Selvam, and R. Babu, "An overview of machine learning and its applications," International Journal of Electrical Sciences & Engineering, vol. 1, pp. 22–24, 2016.

[5] https://www.aretove.com/importance-of-feature-selection-in-machine learning#:~:text=Feature%20selection%20offers%20a%20simple,the%20learning%20model%20or%20data.

[6] Chawla, A. (2022). Phishing website analysis and detection using Machine Learning. International Journal of Intelligent Systems and Applications in Engineering, 10(1), 10–16. https://doi.org/10.18201/ijisae.2022.262

[7] J. Alzubi, A. Nayyar, and A. Kumar, "Machine learning from theory to algorithms: an overview," Journal of Physics: Conference Series, vol. 1142, Article ID 012012, 2018.

[8] Kasani, Payam Hosseinzadeh, et al. "An Evaluation of Machine Learning Classifiers for Prediction of Alzheimer's Disease, Mild Cognitive Impairment and Normal Cognition." 2021 International Conference on Information and Communication Technology Convergence (ICTC). IEEE, 2021.

[9] Bansal, Deepika, et al. "Comparative analysis of various machine learning algorithms for detecting dementia." Procedia computer science 132 (2018): 1497-1502.

[10] Shahbaz, Muhammad, et al. "Classification of Alzheimer's Disease using Machine Learning Techniques." In Data (2019): 296-303.

[11] Leong, Lee Kuok, and Azian Azamimi Abdullah. "Prediction of Alzheimer's disease (AD) using machine learning techniques with boruta algorithm as feature selection method." Journal of Physics: Conference Series. Vol. 1372. No. 1. IOP Publishing, 2019.

[12] Ghazaly, N. M. . (2022). Data Catalogue Approaches, Implementation and Adoption: A Study of Purpose of Data Catalogue. International Journal on Future Revolution in Computer Science &Amp; Communication Engineering, 8(1), 01–04. https://doi.org/10.17762/ijfrcsce.v8i1.2063

[13] Alickovic, Emina, Abdulhamit Subasi, and Alzheimer's Disease Neuroimaging Initiative. "Automatic detection of alzheimer disease based on histogram and random forest." International Conference on Medical and Biological Engineering. Springer, Cham, (2020): 91-96.

[14] Martinez-Murcia, Francisco J., et al. "Studying the manifold structure of Alzheimer's disease: a deep learning approach using convolutional autoencoders." IEEE journal of biomedical and health informatics 24.1 (2020): 17-26.

[15] Alroobaea, Roobaea, et al. "Alzheimer's Disease Early Detection Using Machine Learning Techniques." (2021).

[16] Faouri, S., M. AlBashayreh, and M. Azzeh. "Examining stability of machine learning methods for predicting dementia at early phases of the disease." Decision Science Letters 11.3 (2022): 333-346.

[17] Kavitha, C., et al. "Early-Stage Alzheimer's Disease Prediction Using Machine Learning Models." Frontiers in Public Health 10 (2022).

[18] Khan, Afreen, and Swaleha Zubair. "An improved multi-modal based machine learning approach for the prognosis of Alzheimer's disease." Journal of King Saud University-Computer and Information Sciences (2020).

[19] Marcus, Daniel S., et al. "Open access series of imaging studies: longitudinal MRI data in non-demented and demented older adults." Journal of cognitive neuroscience 22.12 (2010): 2677-2684.

[20] Marcus, Daniel S., et al. "Open Access Series of Imaging Studies (OASIS): cross-sectional MRI data in young, middle aged, non-demented, and demented older adults." Journal of cognitive neuroscience 19.9 (2007): 1498-1507.

[21] Ahmed Cherif Megri, Sameer Hamoush, Ismail Zayd Megri, Yao Yu. (2021). Advanced Manufacturing Online STEM Education Pipeline for Early-College and High School Students. Journal of Online Engineering Education, 12(2), 01–06. Retrieved from http://onlineengineeringeducation.com/index.php/joee/article/view/47

[22] Ogudo, Kingsley A., et al. "A device performance and data analytics concept for smartphones' IoT services and machine-type communication in cellular networks." Symmetry 11.4 (2019): 593.

[23] Geurts, Pierre, Damien Ernst, and Louis Wehenkel. "Extremely randomized trees." Machine learning 63.1 (2006): 3-42.

[24] Natekin, A., and A. Knoll. "Gradient boosting machines, a tutorial, Front. Neurorobot., 7, 21." (2013).

[25] Ying, Cao, et al. "Advance and prospects of AdaBoost algorithm." Acta Automatica Sinica 39.6 (2013): 745-758.

[26] Alhaj, Taqwa Ahmed, et al. "Feature selection using information gain for improved structural-based alert correlation." PloS one 11.11 (2016): e0166017.