

Machine Learning Models for Alzheimer's Disease Prediction: A Comparative Study

R. Arumugam¹ and A. Murugan²

¹Research Scholar, Department of Computer and Information Science, Annamalai University, Annamalai Nagar – 608 002, Tamil Nadu, India

Email: arumugammca848@gmail.com

²Assistant Professor, Department of Computer Science, Periyar Arts College, Cuddalore, (Deputed from Annamalai University, Annamalai Nagar) Tamil Nadu, India

Email: drmuruganapcs@gmail.com

Abstract

Alzheimer's disease (AD) is a progressive neurodegenerative disorder that significantly impacts cognitive and functional abilities. Early detection and accurate prediction of AD progression are critical for effective intervention and management. This study explores the predictive potential of machine learning algorithms, including Linear Regression, Multilayer Perceptron, SMOreg, Random Forest, Random Tree, and REP Tree, applied to Alzheimer's-related datasets. The dataset comprises features such as demographic (ID, gender, handedness, age, education, socioeconomic status), cognitive (MMSE, CDR), and structural (ETIV, NWBV, ASF) attributes. Comprehensive analysis reveals the strengths and limitations of each model in handling the diverse dataset characteristics. The results demonstrate that tree-based methods like Random Forest and REP Tree provide superior accuracy, while neural network-based approaches like Multilayer Perceptron effectively capture nonlinear relationships. This research underscores the importance of integrating cognitive and structural metrics to enhance predictive capabilities, offering valuable insights for early diagnosis and personalized care strategies in Alzheimer's disease.

Keywords- Alzheimer's disease, machine learning, cognitive metrics, brain volume, prediction models, decision trees.

1.0 Introduction and Literature Review

Alzheimer's disease (AD) is a progressive neurodegenerative disorder characterized by the gradual deterioration of memory, cognitive abilities, and daily functioning. It poses a significant challenge globally, affecting millions of individuals and placing a substantial burden on caregivers and healthcare systems. Early detection and accurate prediction of Alzheimer's disease progression are vital for effective management, timely interventions, and personalized care strategies. Machine learning (ML) has emerged as a powerful tool in medical research, providing robust methodologies for analyzing complex datasets and uncovering patterns that traditional statistical methods might overlook.

The predictive modeling of Alzheimer's disease relies on various datasets comprising demographic, cognitive, and structural features. Demographic features, such as age, gender, handedness, education, and socioeconomic status (SES), offer insights into population-specific risk factors. Cognitive assessments,

including the Mini-Mental State Examination (MMSE) and Clinical Dementia Rating (CDR), quantify cognitive decline and dementia severity. Structural metrics derived from brain imaging, such as estimated total intracranial volume (ETIV), normalized whole brain volume (NWBV), and atlas scaling factor (ASF), provide valuable information about brain morphology and atrophy associated with Alzheimer's disease.

This study investigates the application of six machine learning algorithms for the analysis and prediction of Alzheimer's disease progression. The models include Linear Regression, Multilayer Perceptron (MLP), Sequential Minimal Optimization Regression (SMOreg), Random Forest, Random Tree, and Reduced Error Pruning (REP) Tree. Each algorithm offers unique advantages, from capturing linear relationships to modeling complex nonlinear interactions, making them suitable for handling diverse dataset characteristics.

The primary objective of this research is to evaluate the performance of these models in accurately predicting Alzheimer's disease progression and identifying key contributing factors. By integrating cognitive scores and structural metrics, this study aims to enhance the predictive capabilities of machine learning approaches, ultimately contributing to the early detection and effective management of Alzheimer's disease. The findings underscore the importance of leveraging advanced machine learning techniques to address the complexities of neurodegenerative disorders and provide actionable insights for clinical decision-making.

This article serves multiple purposes. Firstly, it outlines the fundamental steps of a CAD system for brain MRI, highlighting the research related to diagnosing brain disorders, especially Alzheimer's disease (AD). The common methods in classification and brain region segmentation are discussed with their respective advantages and disadvantages. Secondly, the article proposes a solution within the realm of multimodal fusion to address the raised issue. It introduces a performance study of multimodal CAD systems based on quantitative measurement parameters, comparing their effectiveness with single MRI modality-based systems. The article underscores advancements in information fusion techniques in medical imaging, highlighting their pros and cons. Finally, it delves into the contributions of multimodal fusion and the significance of hybrid models, summarizing key scientific assertions in the field of brain disease diagnosis [1].

Another study, the author recommends the application of six different machine learning and data mining algorithms, including k-nearest neighbors (k-NN), decision tree (DT), rule induction, Naive Bayes, generalized linear model (GLM), and deep learning, to classify the five different stages of AD using the ADNI dataset. The results show that GLM achieves an accuracy of 88.24% in classifying AD stages, emphasizing the potential of these techniques in early disease detection and diagnosis [2].

Another research effort segmented a dataset into variable sizes and employed three machine learning algorithms: Sequential Minimal Optimization (SMO), Naïve Bayes (NB), and J48, to predict cognitive impairment status. The J48-based model achieved impressive results, with an accuracy of 98.82%, an AUC value of 0.992, and high sensitivity and specificity values, showcasing its effectiveness in classifying Alzheimer's patients [3].

The discussion highlights the focus of existing research on predicting the progression of AD dementia using publicly available datasets encompassing neuroimaging and clinical data. It underscores the utility of clinical data in machine learning-based risk modeling for AD dementia progression, emphasizing the importance of data sharing and reproducibility [4].

A different study, patient data was used to train a model for forecasting disease progression in Mild Cognitive Impairment and Alzheimer's Disease. The study demonstrated the effectiveness of this model in predicting changes in cognitive exam scores and identified specific components predictive of progression [5].

The paper discusses the application of data mining to weather data for golf playing conditions. It presents the results of using seven classification algorithms, with the Random Tree algorithm achieving the highest accuracy of 85.714% [6].

A study employed four machine learning models and feature selection methods to predict the development of AD at an early stage. It concluded that ensemble modeling with selective features offers improved accuracy [7].

The research evaluates various methods for early detection of AD using machine learning techniques, highlighting the challenge of comparing these methods due to variations in data sets and factors such as pre-processing and feature selection [8].

Lastly, the paper proposes a model that involves pre-processing, attribute selection, and classification using association rule mining to distinguish AD from healthy controls, addressing the need for standardized evaluation methods [9].

A different context, the article applies stochastic modeling and data mining to groundwater level, rainfall, population, food grains, and enterprises data to predict groundwater levels accurately [10] and [11].

The study focuses on chronic disease data and assesses the accuracy of five different decision tree algorithms, with the M5P decision tree approach showing superior model-building capabilities [12].

2.0 Backgrounds and Methodologies

A data mining decision tree is a widely used machine learning technique for classification and regression tasks. It visually depicts a sequence of

decisions and their possible outcomes in a tree-like structure. Each internal node represents a decision based on a specific feature, and each branch corresponds to the potential result of that decision. The tree's leaf nodes represent the final decision or the predicted outcome. The "CART" (Classification and Regression Trees) algorithm is the most used algorithm for building decision trees [13].

2.1 Linear Regression

Linear regression is a statistical technique employed to comprehend and forecast the connection between two variables by discovering the optimal straight line that most effectively aligns with the data points. It aids in ascertaining how alterations in one variable correspond to changes in another, proving valuable for predictions and trend recognition.

The core idea of linear regression is to find the best-fitting straight line (also called the "regression line") through a scatterplot of data points. This line represents a linear equation of the form:

$$y = m_x + b$$

... (1)

Where:

- ❖ y is the dependent variable (the one you want to predict or explain).
- ❖ x is the independent variable (the one you're using to make predictions or explanations).
- ❖ m is the slope of the line, representing how much
- ❖ y changes for a unit change in x .

b is the y -intercept, indicating the value of y when x is 0.

2.2 Multilayer Perception

A Multilayer Perceptron (MLP) is an artificial neural network consisting of multiple layers of interconnected nodes or neurons. It's a fundamental architecture in deep learning and is used for various tasks, including classification, regression, and more complex tasks like image recognition and natural language processing. The architecture of an MLP typically includes three types of layers:

- i. **Input Layer:** This layer consists of neurons receiving input data. Each neuron corresponds to a feature in the input data, and the values of these neurons pass through the network.

- ii. **Hidden Layers:** These layers come after the input layer and precede the output layer. They are called "hidden" because their activations are not directly observed in the final output.
- iii. **Output Layer:** This layer produces the network's final output. The number of neurons in the output layer depends on the problem type.

2.3 SMO

SMO stands for "Sequential Minimal Optimization," an algorithm used for training support vector machines (SVMs), machine learning models commonly used for classification and regression tasks. The SMO algorithm is particularly well-suited for solving the quadratic programming optimization problem that arises during the training of SVMs.

- Step 1. **Initialization:** Start with all the data points as potential support vectors and initialize the weights and bias of the SVM.
- Step 2. **Selection of Two Lagrange Multipliers:** In each iteration, the SMO algorithm selects two Lagrange multipliers (associated with the support vectors) to optimize.
- Step 3. **Optimize the Pair of Lagrange Multipliers:** Fix all the Lagrange multipliers except the selected two, and then optimize the pair chosen to satisfy certain constraints while maximizing a specific objective function.
- Step 4. **Update the Model:** After optimizing the selected pair of Lagrange multipliers, update the SVM model's weights and bias based on the new values of the Lagrange multipliers.
- Step 5. **Convergence Checking:** Check for convergence criteria to determine whether the algorithm should terminate.
- Step 6. **Repeat:** If convergence hasn't been reached, repeat steps 2 to 5 until it is.

2.4 Random Forest

Random Forest is a popular machine learning ensemble method for classification and regression tasks. It is an extension of decision trees and is known for its high accuracy, robustness, and ability to handle complex datasets. Random Forest is widely used in various domains, including data science, machine learning, and pattern recognition. The main idea behind Random Forest is to create an ensemble (a collection) of decision trees and combine their predictions to make more accurate and stable predictions. The following steps describe what Random Forest works like.

- ❖ Bootstrap Aggregating (Bagging)
- ❖ Decision Tree Construction
- ❖ Voting for Classification, Averaging for Regression

The key advantages of Random Forest are:

- ❖ Reduced overfitting
- ❖ Robustness
- ❖ Feature Importance

Steps involved in Random Forest

Random Forest is an ensemble learning method combining multiple decision trees to make more accurate and robust predictions for classification and regression tasks. The steps involved in building a Random Forest are as follows:

- Step 1. Data Bootstrapping
- Step 2. Random Feature Subset Selection
- Step 3. Decision Tree Construction
- Step 4. Ensemble of Decision Trees
- Step 5. Out-of-Bag (OOB) Evaluation
- Step 6. Hyperparameter Tuning (optional)

2.5 Random Tree

In machine learning, a Random Tree is a specific type of decision tree variant that introduces randomness during construction. Random Trees are similar to traditional decision trees but differ in how they select the splitting features and thresholds at each node. The primary goal of introducing randomness is to create a more diverse set of decision trees, which can help reduce overfitting and improve the model's generalization performance. Random Trees are commonly used as building blocks in ensemble methods like Random Forests. The critical characteristics of Random Trees are as follows:

- ❖ Random Feature Subset
- ❖ Random Threshold Selection
- ❖ No Pruning
- ❖ Ensemble Methods

Steps involved in Random Tree

- Step 1. Data Bootstrapping:
- Step 2. Random Subset Selection for Features:
- Step 3. Decision Tree Construction:
- Step 4. Voting (Classification) or Averaging (Regression):

2.6 REP Tree

REP (Repeated Incremental Pruning to Produce Error Reduction) Tree is a machine learning algorithm for classification and regression tasks. A decision tree-based algorithm constructs a decision tree using a combination of incremental pruning and error-reduction techniques. The key steps involved in building a REP Tree are as follows:

- ❖ Recursive Binary Splitting
- ❖ Pruning
- ❖ Repeated Pruning and Error Reduction

Steps involved in REP Tree

REP Tree (Repeated Incremental Pruning to Produce an Error Reduction Tree) is a machine learning algorithm for classification and regression tasks. It is an extension of decision trees that incorporates pruning to reduce overfitting and improve the model's generalization performance. Below are the steps involved in building a REP Tree.

- Step 1. Recursive Binary Splitting
- Step 2. Pruning
- Step 3. Repeated Pruning and Error Reduction
- Step 4. Model Evaluation

2.7 Accuracy Metrics

The predictive model's error rate can be evaluated by applying several accuracy metrics in machine learning and statistics. The basic concept of accuracy evaluation in regression analysis is comparing the original target with the predicted one and using metrics like R-squared, MAE, MSE, and RMSE to explain the errors and predictive ability of the model [14]. The R-squared, MSE, MAE, and RMSE are metrics used to evaluate the prediction error rates and model performance in analysis and predictions [15] and [16].

R-squared (Coefficient of determination) represents the coefficient of how well the values fit compared to the original values. The values from 0 to 1 are interpreted as percentages. The higher the value is, the better the model is.

$$R^2 = 1 - \frac{\sum(y_i - \hat{y})^2}{\sum(y_i - \bar{y})^2} \quad \dots (2)$$

MAE (Mean absolute error) represents the difference between the original and predicted values

extracted by averaging the absolute difference over the data set.

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}| \quad \dots (3)$$

RMSE (Root Mean Squared Error) is the error rate by the square root of MSE.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2} \quad \dots (4)$$

Relative Absolute Error (RAE) is a metric used in statistics and data analysis to measure the accuracy of a forecasting or predictive model's predictions. It is particularly useful when dealing with numerical data, such as in regression analysis or time series forecasting.

$$RAE = \frac{\sum |y_i - \hat{y}_i|}{\sum |y_i - \bar{y}|} \quad \dots (5)$$

Root Relative Squared Error (RRSE) is another metric used in statistics and data analysis to evaluate the accuracy of predictive models, especially in the context of regression analysis or time series forecasting.

$$RRSE = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}} \quad \dots (6)$$

Equation 3 to 7 are used to find the model accuracy, which is used to find the model performance and error. Where Y_i represents the individual observed (actual) values, \hat{Y}_i represents the corresponding

individual predicted values, \bar{Y} represents the mean (average) of the observed values and Σ represents the summation symbol, indicating that you should sum the absolute differences for all data points.

3.0 Numerical Illustrations

The corresponding dataset was collected from the open source Kaggle data repository. The alzheimer dataset include 11 parameters which have different categories of data like id, m/f, hand, age, educ, ses, mmse, cdr, etiv, nwbv, asf [17]. A detailed description of the parameters is mentioned in the following Table 1.

ID	Identification
M/F	Gender
Hand	Dominant Hand
Age	Age in years
Educsort	Education Level
SES	Socioeconomic Status
MMSE	Mini Mental State Examination
CDR	Clinical Dementia Rating
eTIV	Estimated Total Intracranial Volume
nWBV	Normalize Whole Brain Volume
ASF	Atlas Scaling Factor

Table 1. Alzheimers sample dataset

ID	M/F	Hand	Age	Educ	SES	MMSE	eTIV	nWBV	ASF	CDR
OAS1_0001_MR1	F	R	74	2	3	29	1344	0.743	1.306	0
OAS1_0002_MR1	F	R	55	4	1	29	1147	0.81	1.531	0
OAS1_0003_MR1	F	R	73	4	3	27	1454	0.708	1.207	0.5
OAS1_0010_MR1	M	R	74	5	2	30	1636	0.689	1.073	0
OAS1_0011_MR1	F	R	52	3	2	30	1321	0.827	1.329	0

Table 2: Machine Learning Models with Correlation coefficient

ML Approaches	Correlation coefficient
Linear Regression	0.9892
Multilayer Perceptron	1.0000

SMOreg	0.9890
Random Forest	0.9871
Random Tree	0.9422
REP Tree	0.9937

Table 3: Machine Learning Models with Mean Absolute Error and Root Mean Squared Error

ML Approaches	MAE	RMSE
Linear Regression	0.0134	0.0189
Multilayer Perceptron	0.0009	0.0013
SMOreg	0.0125	0.0201
Random Forest	0.0165	0.0269
Random Tree	0.0260	0.0431
REP Tree	0.0078	0.0145

Table 4: Machine Learning Models with Relative Absolute Error (%) and Root Relative Squared Error (%)

ML Approaches	RAE (%)	RRSE (%)
Linear Regression	13.1225	14.6484
Multilayer Perceptron	0.8611	0.9771
SMOreg	12.2036	15.5897
Random Forest	16.0507	20.8534
Random Tree	25.4115	33.4519
REP Tree	7.6229	11.2449

Table 5: Machine Learning Models with Time Taken to Build Model (Seconds)

ML Approaches	Time taken (seconds)
Linear Regression	0.5800
Multilayer Perceptron	0.8100
SMOreg	0.2100
Random Forest	0.5600
Random Tree	0.8700
REP Tree	0.0400

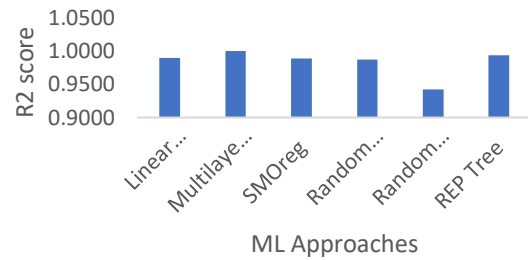


Fig. 1. R2 Score for Machine Learning Approaches

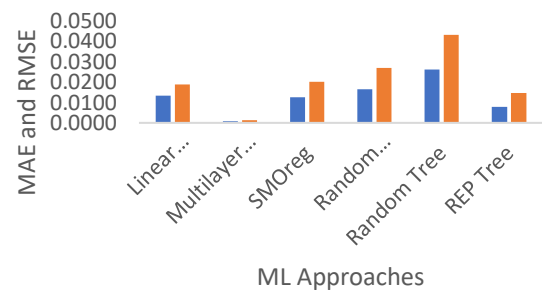


Fig. 2. Machine Learning Models with MAE and RMSE

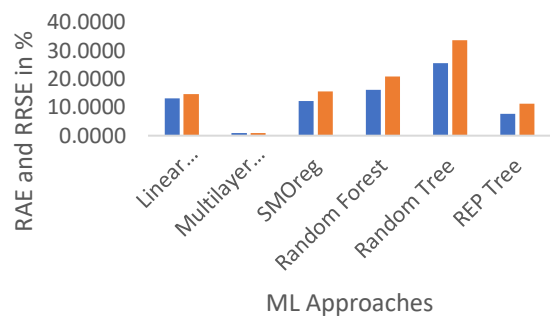


Fig. 3. Machine Learning Models with RAE (%) and RRSE (%)

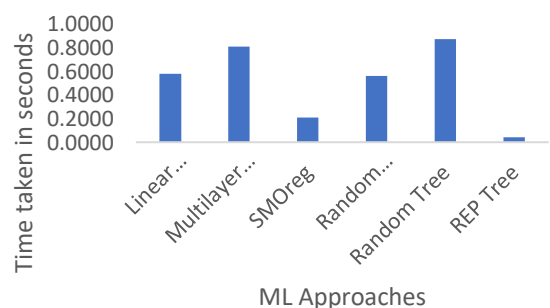


Fig. 4. Machine Learning Models and its Time Taken to Build the Model (Seconds)

4.0 Results and Discussion

In Section 4.0, we delve into the analysis and discussion of the obtained results. Table 1 provides a comprehensive overview of 11 parameters, encompassing various data categories, such as id, gender (m/f), dominant hand, age, education level (educ), socio-economic status (ses), Mini-Mental State Examination (mmse) scores, Clinical Dementia Rating (cdr), Estimated Total Intracranial Volume (etiv), Normalized Whole Brain Volume (nwbv), and Atlas Scaling Factor (asf). These parameters are pivotal in assessing their impact on future predictions. To uncover hidden patterns within the dataset, we employed six distinct machine learning approaches: linear regression, multilayer perceptron, SMOREg, random forest, random tree, and REP tree. These methodologies were instrumental in identifying the most influential parameters for predicting future outcomes. The results, along with numerical illustrations, can be found in Tables 1 to 5 and Figures 1 to 4.

The analyses were based on Equation 2, Table 2, and Figure 1, which facilitated the calculation of the R2 score and correlation coefficients across the 11 parameters. Notably, our numerical illustrations highlight significant variations among these parameters. Particularly, when using the asf parameter, all six approaches demonstrate a strong positive correlation, nearing 0.9.

Furthermore, we assessed model errors using the Mean Absolute Error (MAE) as defined by Equation 3. In this investigation, six distinct machine-learning algorithms were employed, all of which exhibited exceptional error performance, nearing an error rate of 0. Similarly, the Root Mean Square Error (RMSE) was utilized, as described in Equation 4, to measure the disparities between predicted and actual values. In this case, all machine learning approaches yielded outstanding error performance, approaching an error rate of 0. Detailed numerical insights can be found in Table 3 and Figure 2.

Additionally, we employed the Relative Absolute Error (RAE) defined by Equation 5 to assess accuracy by comparing the disparities between predicted and actual values as percentages. This examination encompassed six machine learning classification algorithms, revealing that the random tree approach exhibited the highest error rate, while the other five machine learning approaches consistently delivered

optimal performance with minimal errors. The corresponding error analyses were extended to RRSE, as reflected in Table 4 and Figure 3.

Finally, we examined the time taken by the machine learning approaches. As detailed in Table 5 and Figure 4, Multilayer Perceptron and Random Tree required the most time to address this problem, whereas Random Forest and REP Tree displayed the quickest model development. Linear Regression and the SMOREg approach also exhibited minimal time requirements to advance the model. These findings are consistent with the visual representations provided.

5.0 Conclusion and Future Research

Addresses the constraints within our model, encompassing considerations related to data parameters such as id, gender, hand dominance, age, education level, socio-economic status, cognitive assessment scores, and brain volume metrics, including asf. We also acknowledge potential model-specific factors that may contribute to underperformance and computational constraints that may have influenced model development. In conclusion, we propose potential enhancements and future research directions. These include exploring additional data sources, investigating improved algorithms and hyperparameters, and fine-tuning the model to enhance its overall performance. It is essential to recognize that Alzheimer's research is a complex and ongoing endeavor with no current cure. Nonetheless, our continuous efforts are advancing our understanding of Alzheimer's disease and contributing to the development of therapies and interventions to improve the well-being of individuals affected by this condition.

Reference

1. Lazli, L., Boukadoum, M. and Mohamed, O.A., 2020. A survey on computer-aided diagnosis of brain disorders through MRI based on machine learning and data mining methodologies with an emphasis on Alzheimer disease diagnosis and the contribution of the multimodal fusion. *Applied Sciences*, 10(5), p.1894.
2. Shahbaz, M., Ali, S., Guergachi, A., Niazi, A. and Umer, A., 2019, July. Classification of Alzheimer's Disease using Machine Learning Techniques. In *Data* (pp. 296-303).
3. Hassan, S.A. and Khan, T., 2017. A machine learning model to predict the onset of Alzheimer disease using potential cerebrospinal

- fluid (CSF) biomarkers. *International Journal of Advanced Computer Science and Applications*, 8(12).
4. Kumar, S., Oh, I., Schindler, S., Lai, A.M., Payne, P.R. and Gupta, A., 2021. Machine learning for modeling the progression of Alzheimer disease dementia using clinical data: a systematic literature review. *JAMIA open*, 4(3), p.ooab052.
 5. Fisher, C.K., Smith, A.M. and Walsh, J.R., 2019. Machine learning for comprehensive forecasting of Alzheimer's Disease progression. *Scientific reports*, 9(1), pp.1-14.
 6. Rajesh, P. and Karthikeyan, M., 2017. A comparative study of data mining algorithms for decision tree approaches using the Weka tool. *Advances in Natural and Applied Sciences*, 11(9), pp.230-243.
 7. Battineni, G., Chintalapudi, N., Amenta, F. and Traini, E., 2020. A comprehensive machine-learning model applied to magnetic resonance imaging (mri) to predict alzheimer's disease (ad) in older subjects. *Journal of Clinical Medicine*, 9(7), p.2146.
 8. Ammar, R.B. and Ayed, Y.B., 2018, October. Speech processing for early Alzheimer disease diagnosis: machine learning based approach. In 2018 IEEE/ACS 15th International Conference on Computer Systems and Applications (AICCSA) (pp. 1-8). IEEE.
 9. Khan, A. and Usman, M., 2015, November. Early diagnosis of Alzheimer's disease using machine learning techniques: A review paper. In 2015 7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management (IC3K) (Vol. 1, pp. 380-387). IEEE.
 10. Rajesh, P., Karthikeyan, M. and Arulpavai, R., 2019, December. Data mining approaches to predict the factors that affect the groundwater level using a stochastic model. In AIP Conference Proceedings (Vol. 2177, No. 1). AIP Publishing.
 11. Rajesh, P. and Karthikeyan, M., 2019. Data mining approaches to predict the factors that affect agriculture growth using stochastic models. *International Journal of Computer Sciences and Engineering*, 7(4), pp.18-23.
 12. Rajesh, P., Karthikeyan, M., Santhosh Kumar, B. and Mohamed Parvees, M.Y., 2019. Comparative study of decision tree approaches in data mining using chronic disease indicators (CDI) data. *Journal of Computational and Theoretical Nanoscience*, 16(4), pp.1472-1477.
 13. Kohavi, R., & Sahami, M. (1996). Error-based pruning of decision trees. In *International Conference on Machine Learning* (pp. 278-286).
 14. Akusok, A. (2020). What is Mean Absolute Error (MAE)? Retrieved from <https://machinelearningmastery.com/mean-absolute-error-mae-for-machine-learning/>
 15. S. M. Hosseini, S. M. Hosseini, and M. R. Mehrabian, "Root mean square error (RMSE): A comprehensive review," *International Journal of Applied Mathematics and Statistics*, vol. 59, no. 1, pp. 42–49, 2019.
 16. Chi, W. (2020). Relative Absolute Error (RAE) – Definition and Examples. Medium. <https://medium.com/@wchi/relative-absolute-error-rae-definition-and-examples-e37a24c1b566>
 17. <https://www.kaggle.com/code/hyunseokc/detecting-early-alzheimer-s-input>