

Machine Learning Approaches for Prediction of Alzheimer's Disease

Kadriye Filiz BALBAL ^{1*}

Abstract

Alzheimer's Disease (AD) is a disorder that significantly impacts an individual's behavior, memory, and cognitive functions, ultimately leading to a loss of independence. Early and accurate diagnosis of AD is critical to mitigating its progression and improving patient outcomes, especially as no definitive cure is currently available. This study investigates the application of machine learning algorithms to predict and diagnose AD based on patient symptoms and clinical data. The dataset used in this research includes comprehensive health information from 2,149 patients, with 35 features covering demographic, lifestyle, and medical factors, and no missing values. Seven widely recognized machine learning algorithms—KNN, GNB, SVM, DT, RF, AdaBoost, and XGBoost—were evaluated to determine their effectiveness in disease prediction. Performance was assessed using recall, precision, accuracy, and F1-score metrics, providing a robust evaluation of each model. XGBoost achieved the highest accuracy rate of 95.35%, highlighting its superior predictive capability, while KNN recorded the lowest accuracy at 75.54%. The results demonstrate the strength of machine learning algorithms, particularly ensemble methods like XGBoost, in analyzing complex clinical data for the early detection of Alzheimer's Disease. These findings underscore the critical role of machine learning in enhancing diagnostic accuracy and enabling timely interventions, which are essential for improving the quality of life for individuals at risk of Alzheimer's Disease.

Keywords: *Alzheimer's disease; classification; disease prediction; machine learning.*

1. Introduction

Alzheimer's Disease (AD) is a progressive neurodegenerative disorder that often manifests with symptoms resembling the natural effects of aging, making it challenging to distinguish in its early stages. Symptoms such as forgetfulness, confusion, or behavioral changes like stress and paranoia can overlap with other conditions or be dismissed as normal aging processes. Primarily affecting individuals aged 65 and older, AD progressively deteriorates cognitive functions, eventually impairing the ability to perform daily activities independently [1]. The disease profoundly impacts memory, thinking, and reasoning abilities, gradually diminishing the quality of life for both patients and their caregivers. Despite extensive research, there is currently no definitive cure for AD, and it remains a significant global health challenge. However, evidence suggests that early and accurate diagnosis, coupled with timely interventions, can slow the disease's progression and alleviate its symptoms [2]. By delaying the onset of severe cognitive decline, early diagnosis not only enhances patient outcomes but also reduces the emotional and financial burden on families and healthcare systems. Therefore, identifying AD in its early stages through advanced diagnostic techniques is of paramount importance, as it enables effective management strategies and improves the quality of life for affected individuals [3].

In this study, the prediction of Alzheimer's Disease (AD) status was conducted using machine learning algorithms that are widely recognized in the literature for their effectiveness in disease diagnosis and prediction. The analyses aimed to determine the presence of AD by utilizing various data types, including demographic information, medical history, cognitive and functional evaluations, lifestyle factors, clinical measurements, and patient-reported symptoms. The structure of the paper is organized as follows: Section 2 provides a review of related work, Section 3 details the dataset description and the methodology employed, Section 4 presents the results of the analyses, and Section 5 concludes the study with a discussion of the findings and their implications.

2. Related Works

The role of machine learning (ML) in disease prediction, early diagnosis, and prevention has been extensively emphasized in the literature, highlighting its potential to revolutionize healthcare by improving diagnostic accuracy and enabling timely interventions [4-10]. Numerous studies have explored the application of ML techniques to analyze complex medical datasets, demonstrating their versatility and effectiveness. For example, in [5], where the critical importance of early diagnosis in Alzheimer's Disease (AD) was underscored, predictions were performed using DT and GNB models, showcasing their ability to identify patterns associated with AD.

*Corresponding author

Similarly, in [6], an audio dataset from the UCI Machine Learning Repository was employed to predict diseases using four different ML algorithms. This study demonstrated the adaptability of ML techniques to various data modalities, such as audio features, further expanding their applicability in medical diagnostics.

In [7], a broader investigation was undertaken by predicting diseases such as diabetes, breast cancer, and heart disease using DT and GNB models. The study aimed to identify the most effective algorithm for disease prediction through a comparative analysis, emphasizing the cross-domain applicability of these methods in addressing diverse medical conditions. Furthermore, [8] utilized two supervised ML algorithms for disease prediction, achieving 87% accuracy with GNB and 91% with DT, further substantiating the capability of these methods in clinical data analysis.

In [9], a comprehensive evaluation was conducted using six classification algorithms: Multi-Layer Perceptron (MLP), Logistic Regression (LR), Extremely Randomized Trees Classifier (ERT), SVM, RF, and Gradient Boosting Classifier (GBC), to classify heart disease. This study, performed on the Cleveland dataset containing 14 features, revealed that MLP and SVM achieved the highest performance with an accuracy rate of 91.7%. These findings highlighted the potential of advanced ML algorithms in achieving high precision and reliability in disease classification.

Additionally, in [10], ML algorithms were applied to hospital data collected in Andhra Pradesh, India, between 2018 and 2020, to evaluate their effectiveness in disease prediction. Among the algorithms tested, AdaBoost and KNN were identified as the most successful, demonstrating their ability to extract valuable insights from real-world clinical datasets. Collectively, these studies underscore the importance of ML in disease prediction, highlighting the role of algorithm selection, dataset characteristics, and feature engineering in optimizing model performance. Such findings reveal the transformative potential of ML techniques in healthcare, particularly in addressing diagnostic challenges and enhancing early detection strategies.

3. Data and Methodology

3.1. Dataset

The Alzheimer's Disease dataset was obtained from Kaggle [11]. The dataset with no missing values contains a total of 2,149 rows and 35 columns. The columns in the dataset include information on individuals' demographic characteristics, lifestyle factors, health history, and cognitive status. 16 of the variables are categorical (such as gender, smoking, Alzheimer's diagnosis), and 19 are numerical (such as age, BMI, cholesterol levels). While categorical variables allow patients to be divided into groups, numerical variables allow for more detailed analyses. The dataset provides a wide range of data to understand the risk factors and symptoms of Alzheimer's disease. This diversity provides a solid basis for both descriptive analyses and more advanced statistical methods. Detailed information about the dataset is provided in Table 1.

Table 1. Alzheimer's Disease Dataset Description

Category	Subcategory	Description
Patient Identification	PatientID	Unique identifier assigned to each patient (4751-6900).
Demographic Details	Age	Patients' ages range between 60 and 90 years.
	Gender	0: Male, 1: Female.
	Ethnicity	0: Caucasian, 1: African-American, 2: Asian, 3: Other.
	Education Level	0: None, 1: High School, 2: Undergraduate, 3: Higher.
Lifestyle Factors	BMI	Body Mass Index ranging from 15 to 40.
	Smoking	0: No, 1: Yes.
	AlcoholConsumption	Weekly alcohol consumption ranging from 0 to 20 units.
	PhysicalActivity	Weekly physical activity in hours, ranging from 0 to 10.
	DietQuality	Diet quality score ranging from 0 to 10.
Medical History	SleepQuality	Sleep quality score ranging from 4 to 10.
	FamilyHistoryAlzheimer	0: No, 1: Yes (family history of Alzheimer's disease).
	Cardiovascular Disease	0: No, 1: Yes.
	Diabetes	0: No, 1: Yes.

	Depression Head Injury Hypertension	0: No, 1: Yes. 0: No, 1: Yes (history of head injury). 0: No, 1: Yes.
Clinical Measurements	SystolicBP DiastolicBP TotalCholesterol LDLCholesterol HDLCholesterol Triglycerides	Systolic blood pressure ranging from 90 to 180 mmHg. Diastolic blood pressure ranging from 60 to 120 mmHg. Total cholesterol levels ranging from 150 to 300 mg/dL. Low-density lipoprotein cholesterol levels ranging from 50 to 200 mg/dL. High-density lipoprotein cholesterol levels ranging from 20 to 100 mg/dL. Triglyceride levels ranging from 50 to 400 mg/dL.
Cognitive and Functional Assessments	MMSE FunctionalAssessment MemoryComplaints BehavioralIssues ADL	Mini-Mental State Examination score ranging from 0 to 30 (lower scores indicate cognitive impairment). Functional assessment score ranging from 0 to 10 (lower scores indicate greater impairment). 0: No, 1: Yes (presence of memory complaints). 0: No, 1: Yes (presence of behavioral issues). Activities of Daily Living score ranging from 0 to 10 (lower scores indicate greater impairment).
Symptoms	Confusion Disorientation PersonalityChanges DifficultyCompletingTasks Forgetfulness	0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes. 0: No, 1: Yes.
Diagnosis Information	Diagnosis	0: No, 1: Yes (Alzheimer's diagnosis).
Confidential Information	DoctorInCharge	Information about the doctor in charge (specified as 'XXXConfid' for all patients).

Before proceeding to the analysis, the 'PatientID', and 'DoctorInCharge' columns, which were not important for the analysis, were deleted. The 'Diagnosis Information' column was determined as the target variable and the analysis continued with 32 columns. The distribution of the demographic features Gender, Ethnicity, and Educational Level in the dataset is given in Figure 1.

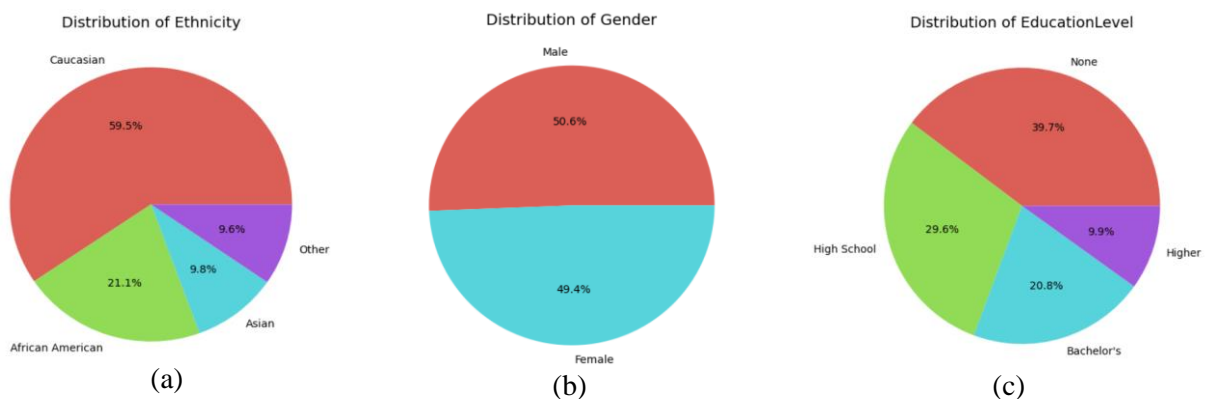


Figure 1. *Distribution of Demographic Features.*

According to Figure 1, 50.6% of the participants in the dataset are male and 49.4% are female. In terms of ethnicity, most of the participants (59.5%) are Caucasian, 21.1% are African American, 9.8% are Asian, and 9.6% are Other. In the dataset, there are 29.6% with High School education level and 20.8% with Bachelor's education level. In addition, 9.9% have Higher education level and 39.7% have neither of these education levels.

3.2. Methodology

This section introduces the machine learning algorithms employed for predicting AD. This study utilized seven widely recognized machine learning algorithms: DT, SVM, GNB, KNN, RF, AdaBoost, and XGBoost. Each algorithm's details and characteristics are outlined below.

K-Nearest Neighbors (KNN) is a supervised learning algorithm commonly applied to regression and classification tasks. The algorithm's performance is governed by the parameter k , which specifies the number of neighbors considered in predictions. While KNN is advantageous for its simplicity and low computational requirements, its performance deteriorates when applied to large datasets due to increased computational costs [12, 13].

Gaussian Naïve Bayes (GNB) operates on the assumption that the features within a class are distributed according to a Gaussian distribution. During training, the algorithm calculates the mean and standard deviation for each class and uses these parameters to estimate probabilities for continuous variables [14].

Support Vector Machine (SVM) is a robust algorithm particularly effective for classification problems in high-dimensional datasets. It supports both linear and nonlinear classification through the use of kernel functions, enabling adaptability to various data structures. However, SVM is computationally intensive and requires complex mathematical operations [15].

Decision Tree (DT) is a tree-based algorithm capable of handling both numerical and categorical data. It excels in capturing intricate interactions between variables and demonstrates robustness against certain levels of noise and inconsistencies in the data [16].

Random Forest (RF) is a powerful ensemble learning algorithm that constructs multiple decision trees to mitigate the high variance observed in single-tree models. It is particularly effective for handling high-dimensional datasets, providing stable and reliable classification performance [17].

Adaptive Boosting (AdaBoost) is an ensemble technique designed to improve the classification performance of weak learners, particularly in scenarios with imbalanced datasets. Its adaptive nature iteratively adjusts the weights of misclassified instances to enhance overall accuracy [18].

Extreme Gradient Boosting (XGBoost) is a high-performance gradient boosting library widely recognized for its efficiency and scalability. XGBoost is particularly effective in classification and prediction problems due to its fast processing capabilities, extensibility, and robust generalization properties [19].

The implementation of these machine learning algorithms for AD prediction was carried out using Python in the Google Colab environment. The dataset, obtained from the Kaggle platform in CSV format, was processed and analyzed using the pandas and numpy libraries. Performance evaluation metrics were imported from the sklearn library, while graphical visualizations were generated using the matplotlib and seaborn libraries. These tools provided a comprehensive framework for data preprocessing, algorithm implementation, and result interpretation in the study.

3.3. Metrics

To evaluate the performance of machine learning algorithms, recall, precision, F1 score, and accuracy metrics were used. These metrics are four evaluation criteria that are widely used in prediction and classification problems [20]. As seen in Equations (1)-(4), the metrics are calculated depending on the values of TN (true negative), TP (true positive), FN (false negative), and FP (false positive).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F1 \text{ score} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

4. Results

In this part of the study, the results obtained from the ML algorithms applied to the AD prediction problem are evaluated in terms of different metrics and presented. The accuracy results obtained from seven different machine learning algorithms are presented and compared in Table 2.

Table 2. Comparison of ML methods in accuracy.

Model	Accuracy
RF	94.12
DT	91.33
SVM	83.59
KNN	75.54
AdaBoost	89.78
GNB	76.47
XGBoost	95.35

In Table 2, the accuracy rates of the machine learning algorithms applied in the study are compared. XGBoost stood out as the most successful model on this dataset, exhibiting the highest performance with an accuracy rate of 95.35%. While the RF algorithm closely follows XGBoost with an accuracy rate of 94.12%, AdaBoost showed a moderate performance with 89.78%. While SVM showed a relatively lower success with an accuracy rate of 83.59%, KNN and GNB showed the lowest performances with accuracy rates of 75.54% and 76.47%, respectively. The results show that the ensemble methods XGBoost and RF are strong options for more complex models.

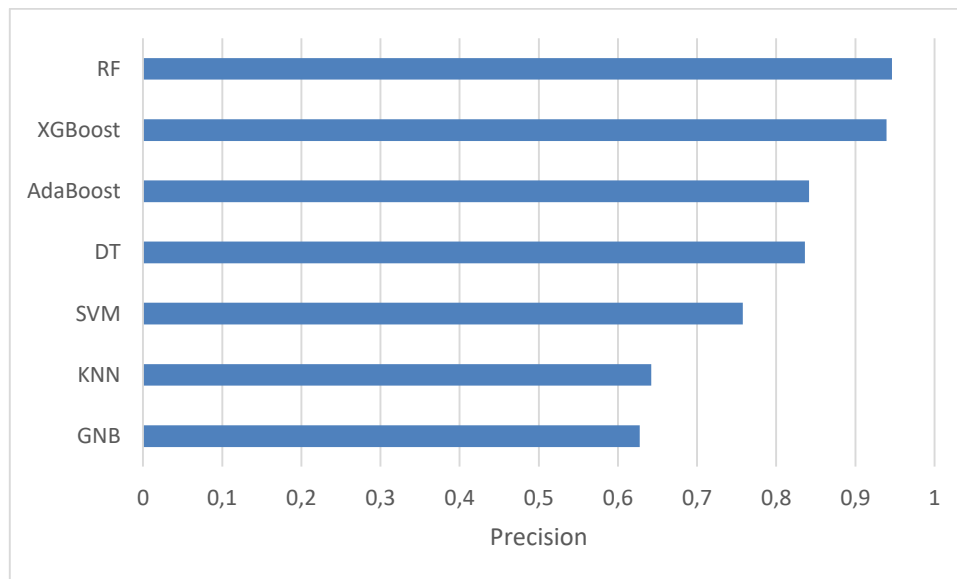


Figure 2. Comparison of ML methods in precision.

Figure 2 compares the precision values of different machine learning algorithms implemented in this study. Precision measures the success of a model in minimizing false positive results by expressing the ratio of true positive predictions to total positive predictions. When Figure 2 is examined, it is seen that XGBoost and RF algorithms reach high precision values and perform better than other models. This shows that these two algorithms can predict true positive results with high accuracy in Alzheimer's diagnosis. In contrast, it is understood that the precision values of KNN and GNB algorithms are relatively lower. AdaBoost and SVM provide balanced results in terms of precision, exhibiting a moderate performance. In general, these results show that XGBoost and RF are effective not only in accuracy but also in reducing false positive predictions. These results emphasize the importance of the precision metric during model selection, especially in areas such as medical diagnosis, where false positives are critical.

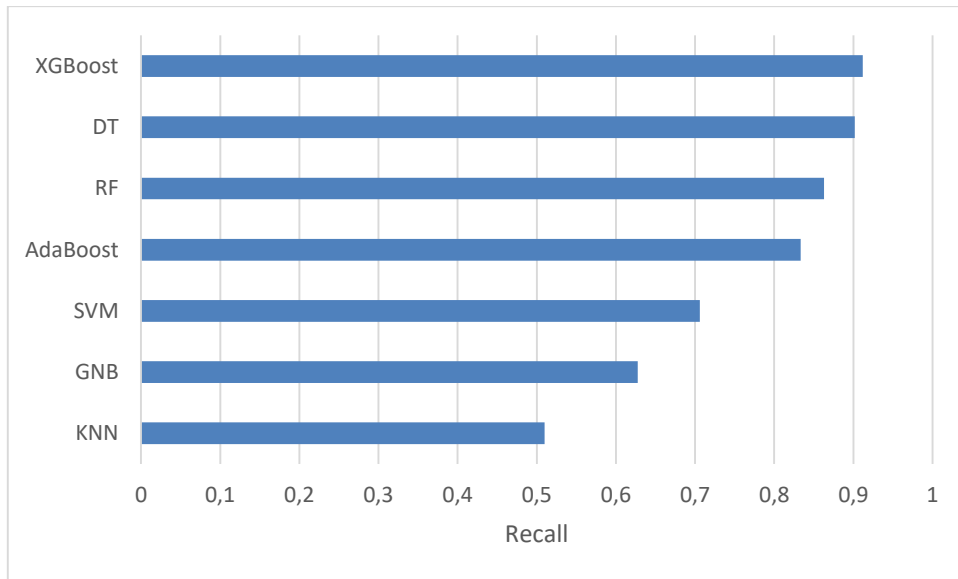


Figure 3. Comparison of ML methods in recall.

Figure 3 compares the recall values of different machine learning algorithms implemented in this study. Recall measures the rate at which a model correctly predicts all true positives and is important for minimizing false negatives (missed positives). As can be seen from the graph, XGBoost, DT, and RF algorithms show superior performance in capturing true positives with high recall values. This highlights the ability of these models to correctly identify patients in a critical diagnosis such as Alzheimer's. AdaBoost and SVM, which showed moderate performance, showed relatively good recall values, but fell behind the most successful models in this metric. KNN and GNB, on the other hand, have lower recall values, indicating that these models missed some of the positive cases. In general, high recall values are of great importance, especially in cases such as disease diagnosis, where missing positive classes is critical. These results show that XGBoost and RF are reliable models, showing a balanced performance not only in terms of accuracy and precision, but also in terms of recall.

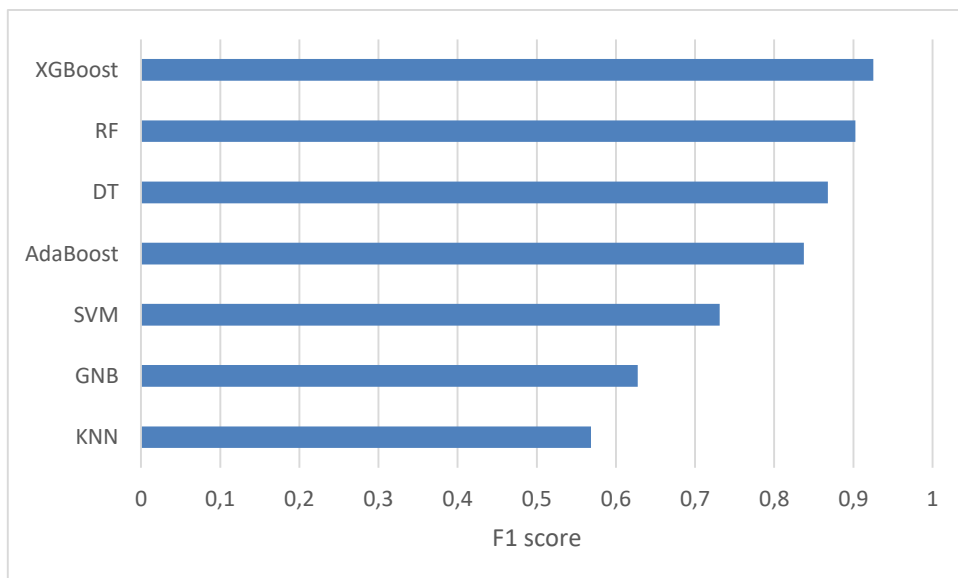


Figure 4. Comparison of ML methods in F1 score.

Figure 4 compares the F1-score values of different machine learning algorithms implemented in this study. Since F1-score is the harmonic mean of precision and recall metrics, it measures the success of a model in correctly predicting the positive class and its ability to deal with imbalanced classes in this process in a holistic way. According to the graphical results, XGBoost and RF achieved the highest F1-score values, which shows that both algorithms effectively detect positive classes while minimizing false positive predictions. It is clear that these models successfully maintain class balance and their overall generalization performance is superior. Although DT and AdaBoost performed reasonably in terms of F1-score, they fell behind the best performing models. On the

other hand, the low F1-score values of models such as KNN and GNB show that these algorithms cannot distinguish between imbalanced classes effectively enough. SVM achieved a moderate success and was relatively successful in finding a balance between precision and recall. This analysis reveals that F1-score is a critical metric and provides a comprehensive evaluation when evaluating model performance, especially on datasets where class imbalance is significant.

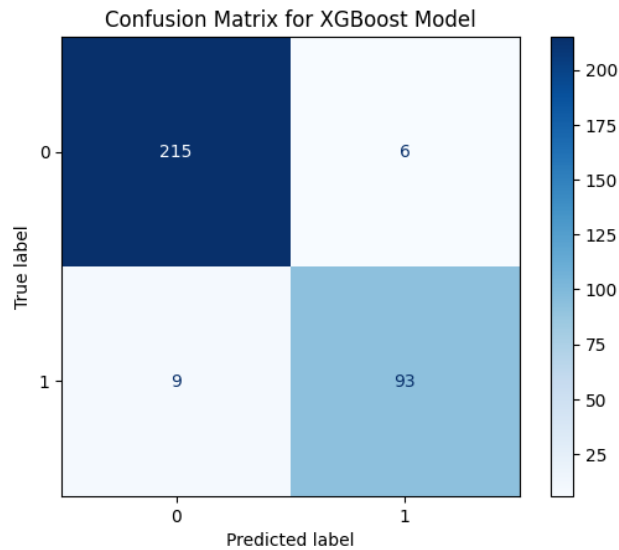


Figure 5. Confusion matrix for XGBoost model.

The confusion matrix in Figure 5 demonstrates the performance of the classification model made with the XGBoost algorithm. According to the matrix, the model predicted 215 true negatives and 93 true positives. However, there were 6 false positives and 9 false negatives among the incorrect predictions. This displays that the model is quite successful in correctly predicting the negative class, but makes relatively more errors for the positive class. When evaluated in general, it is clearly seen that the accuracy rate and prediction performance of the model are high.

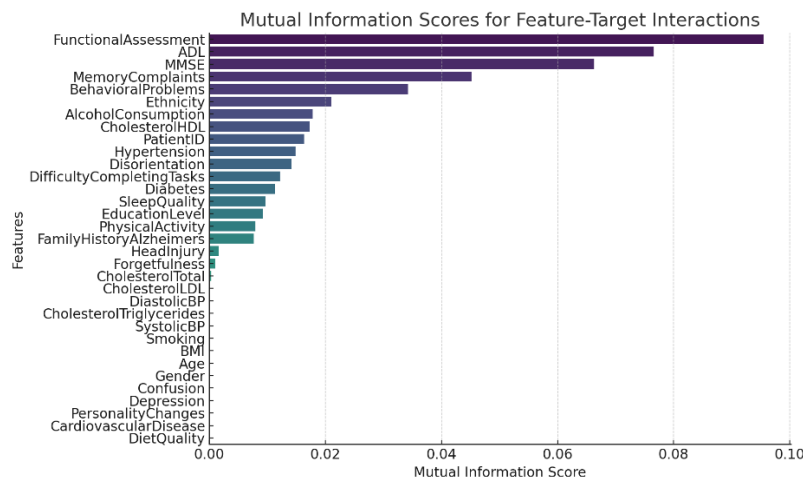


Figure 6. Mutual Information Scores For Feature-Target Interactions.

The graph in Figure 6 evaluates the relationships between the target variable (Alzheimer's diagnosis) and the features using Mutual Information Score. As seen in Figure 6, FunctionalAssessment, ADL (Activities of Daily Living) and MMSE (Mini Mental State Assessment) scores are the features with the strongest relationship with the target variable. This situation reveals that cognitive and functional assessments are critical indicators for AD. Less effective features include factors such as CholesterolTotal, Smoking, BMI, and Age, which indicates that these variables provide limited information in diagnosing the disease. According to the Mutual Information Score results, cognitive status and daily life functions are among the most important factors for both risk assessment and early diagnosis.

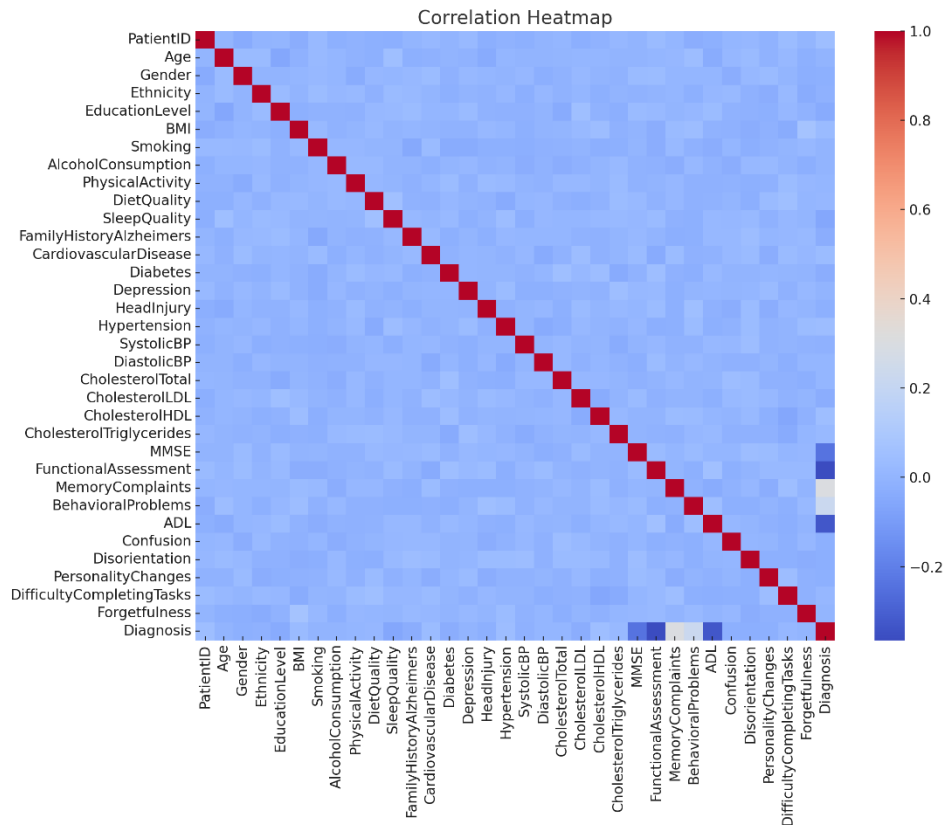


Figure 7. Correlation Heatmap

The correlation heatmap in Figure 7 visualizes the relationships between numerical variables in the dataset. The graph represents the positive or negative correlation of each variable with other variables through color intensity. Specifically, red tones represent high positive correlation (close to 1), while blue tones represent negative correlation (close to -1). Overall, a heatmap is an important tool for understanding which variables are related to each other and for taking these relationships into account during modeling.

5. Conclusion

The findings obtained from this study indicate that machine learning algorithms have demonstrated significant success in predicting Alzheimer's Disease. The experiments were conducted using a dataset that is publicly available on the Kaggle platform, which includes comprehensive demographic and clinical features relevant to Alzheimer's diagnosis. This dataset consists of 2,149 patient records and 35 features, with no missing values, ensuring the reliability and completeness of the data for the analysis. By utilizing this dataset, the study aimed to assess the predictive performance of widely recognized machine learning algorithms frequently employed in disease prediction and classification studies in the literature.

The machine learning algorithms evaluated in this study include XGBoost, KNN, GNB, DT, RF, AdaBoost, and SVM. Each algorithm was assessed using four essential performance metrics: accuracy, precision, recall, and F1-score. The XGBoost algorithm achieved the highest accuracy, with a rate of 95.35%, demonstrating its superior ability to predict Alzheimer's Disease compared to other methods. This finding aligns with previous research, where ensemble learning methods like XGBoost have shown notable effectiveness in handling complex and multidimensional data [8, 9, 21]. In contrast, the lowest accuracy was observed in the KNN algorithm, with a rate of 75.54%. Although KNN performed relatively less effectively than the other models, it is notable that all algorithms achieved accuracy rates above a certain threshold, highlighting their overall competence in predicting Alzheimer's Disease.

The results of this study emphasize the potential of ML algorithms as powerful tools for the early diagnosis and prediction of AD. The superior performance of XGBoost, in particular, emphasizes the importance of using advanced ensemble methods to capture complex patterns and interactions in clinical data. Furthermore, the inclusion of diverse demographic and clinical features in the dataset highlights the significance of integrating multi-dimensional data in predictive modeling. These findings suggest that machine learning-driven predictions

can play a pivotal role in facilitating early detection and guiding timely intervention strategies, ultimately contributing to improved patient outcomes. The study not only confirms the viability of machine learning algorithms in medical diagnosis but also provides valuable insights into their application in addressing critical challenges in healthcare.

References

- [1] D.M. Khan, N. Yahya, N. Kamel, I. Faye, "Automated diagnosis of major depressive disorder using brain effective connectivity and 3D convolutional neural network," *IEEE Access*, 9, pp. 8835-8846, 2021, 10.1109/ACCESS.2021.3049427
- [2] M. Sudharsan and G. Thailambal. "Alzheimer's disease prediction using machine learning techniques and principal component analysis (PCA)," *Materials Today: Proceedings* 81, pp. 182-190, 2023.
- [3] A. Association, "2019 Alzheimer's disease facts and figures", *Alzheimer's & Dementia*, 15 (3), pp. 321-387, 2019.
- [4] C.K. Gomathy and A. Rohith Naidu, "The Prediction Of Disease Using Machine Learning Techniques", *International Journal of Scientific Research in Engineering and Management (IJSREM)*, vol. 5, no. 7, 2021.
- [5] C. R. Mallela, L. R. Bhavani and B. Ankayarkanni, *Disease Prediction Using Machine Learning Techniques*, IEEE, pp. 962-966, 2021.
- [6] T.V. Sriram, M.V. Rao, G.S. Narayana, D.S.V.G. Kaladhar and T.P.R. Vital, "Intelligent Parkinson Disease Prediction Using Machine Learning Algorithms", *International Journal of Engineering and Innovative Technology (IJEIT)*, vol. 3, no. 3, September 2013.
- [7] K. Gomathi and D. Shanmuga Priya, *Multi Disease Prediction using Data Mining Techniques*, 2016.
- [8] A. Gavhane, G. Kokkula, I. Pandya and K. Devatkar, *Prediction of Heart Disease Using Machine Learning Algorithms*, 2018.
- [9] S. Arunachalam, "Cardiovascular Disease Prediction Model using Machine Learning Algorithms", *International Journal for Research in Applied Science & Engineering Technology*, vol. 8, no. VI, June 2020, ISSN 2321-9653.
- [10] A.D., Praveen, T.P., Vital, D., Jayaram and L.V. Satyanarayana, "Intelligent Liver Disease Prediction (ILDLP) System Using Machine Learning Models". *Intelligent Computing in Control and Communication. Lecture Notes in Electrical Engineering*, vol 702, 2021. Springer, Singapore. https://doi.org/10.1007/978-981-15-8439-8_50.
- [11] R.E. Kharoua, "Alzheimer's Disease Dataset", 2024. Kaggle. <https://doi.org/10.34740/KAGGLE/DSV/8668279>.
- [12] K. Taunk, S. De, S. Verma and A. Swetapadma, "A Brief Review of Nearest Neighbor Algorithm for Learning and Classification," 2019 International Conference on Intelligent Computing and Control Systems (ICCS), Madurai, India, 2019, pp. 1255-1260, doi: 10.1109/ICCS45141.2019.9065747.
- [13] I. Triguero, D. García-Gil, J. Maillo, J. Luengo, S. García and F. Herrera, "Transforming big data into smart data: An insight on the use of the k-nearest neighbors algorithm to obtain quality data", *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(2), e1289, 2019.
- [14] H. Kamel, A. Dhahir and J.M. Al-Tuwajjari. "Cancer classification using gaussian naive bayes algorithm." 2019 international engineering conference (IEC). IEEE, 2019.
- [15] S. Suthaharan and S. Shan, "Support vector machine." *Machine learning models and algorithms for big data classification: thinking with examples for effective learning*, 207-235, 2016.
- [16] K. Nong, H. Zhang and Z. Liu, "Comparative Study of Different Machine Learning Models for Heat Transfer Performance Prediction of Evaporators in Modular Refrigerated Display Cabinets", *Energies*, 17, 6189, 2024. <https://doi.org/10.3390/en17236189>
- [17] A. T. Azar, H. I. Elshazly, A. E. Hassanien and A. M. Elkorany, "A random forest classifier for lymph diseases." *Computer methods and programs in biomedicine*, 113(2), 465-473, 2014.
- [18] W. Wang and S. Dongchu, "The improved AdaBoost algorithms for imbalanced data classification." *Information Sciences* 563, 358-374, 2021.
- [19] S. Li and Z. Xiaojing, "Research on orthopedic auxiliary classification and prediction model based on XGBoost algorithm." *Neural Computing and Applications* 32.7 ,1971-1979, 2020.
- [20] P. Iacobescu, V. Marina, C. Anghel, and A-D. Anghel, "Evaluating Binary Classifiers for Cardiovascular Disease Prediction: Enhancing Early Diagnostic Capabilities. *Journal of Cardiovascular Development and Disease*", 11(12):396, 2024. <https://doi.org/10.3390/jcdd11120396>.
- [21] P. Pranjali, S. Mallick, A. Das, A. Negi and M.R. Panda, "Alzheimer's Disease Prediction Using Modern Machine Learning Techniques." 2024 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC). IEEE, 2024.