# Elucidating Encephalitic Disintegration Morbidity Speculations Through Machine Learning

**Mrs.S.Nithiyabharathi**
*Assistant Professor*
*Department of CSE*
*IFET college of engineering*
*Villupuram, India*

**T.Dharshini Vasagan**
*Student*
*Department of CSE*
*IFET college of engineering*
*Villupuram, India*

 *Abstract*: The lifelong neurodevelopmental disorder known as Encephalitic Disintegration Morbidity Speculations (EDMS) usually shows up in early childhood and presents as behavioral, linguistic, and social difficulties. Advanced machine learning approaches are used in this study, and AdaBoost continuously performs better in improving the accuracy of EDMS prediction. The best prediction strategy that works for a wide range of age groups is created when Principal Component Analysis (PCA) and AdaBoost are combined. This research emphasizes the need for early EDMS detection by focusing on the crucial 2–4-month period after a kid is born. In early EDMS diagnosis, especially in children, the suggested ensemble-based model consistently performs better than baseline machine learning techniques, exhibiting superior diagnostic accuracy, precision, recall, and F1-Score. These results show promise for early interventions and improved outcomes for affected individuals and their families, and they represent a major advancement in the improvement of EDMS diagnostic tools. This research provides promise for reducing the long-term effects of this complex neurodevelopmental disorder and improving the quality of life for individuals with EDMS by supporting early detection and intervention techniques.

*IndexTerms* - Encephalitic Disintegration Morbidity Speculations (EDMS), neurodevelopmental disorder, AdaBoost, data imbalance, Ensemble–based model, Principal Component Analysis (PCA).

## INTRODUCTION

Encephalitic Disintegration Morbidity Speculations (EDMS) is a neurodevelopmental illness that affects people's behaviour and communication considerably. People with EDMS have a unique diversity of cognitive capacities. Around one in 160 children worldwide are estimated to have EDMS, but unfortunately, those who have the condition frequently face prejudice and social stigma. To improve early detection and intervention, computer scientists and healthcare professionals have worked together to design screening procedures and apply cutting-edge machine learning algorithms. Specifically, this research focuses on integrating feature selection strategies, data balancing approaches such as the Synthetic Minority Over-sampling Technique (SMOTE), and machine learning classifiers to create a state-of-the-art EDMS detection model. The ultimate objective is to improve the quality of life for people impacted by this complex neurodevelopmental disorder by deepening our understanding of EDMS and utilizing machine learning's potential for early detection.

### 1.1.1 The Intricacy of EDMS

Encephalitic Disintegration Morbidity Speculations (EDMS) are complicated in many ways, beginning with their extensive cognitive capacities. Some people with EDMS have IQs below 70 and mental deficiencies, some have IQs between 71 and 85, and a significant percentage have IQs between average and above average. This wide range of cognitive abilities emphasizes the complex and varied nature of this neurodevelopmental disorder, highlighting the fact that no two people with EDMS have the same strengths and limitations in their cognitive abilities. A wide range of obstacles confront people with EDMS, including as learning disabilities that impede their ability to succeed academically, mental health conditions like anxiety and despair, and physical restrictions that limit their ability to interact with the outside world. The WHO estimates that EDMS affects one in 160 children worldwide. This highlights the critical need for thorough study, creative interventions, and international collaboration to comprehend this intricate neurodevelopmental condition and improve the lives of those impacted.

**1.1.2 Research Objectives**

This work intends to improve the use of modern machine learning approaches for the early detection of Encephalitic Disintegration Morbidity Speculations (EDMS), particularly in pediatric cases. To achieve accurate prediction, the research will compare machine learning methods in detail and use the Synthetic Minority Over-Sampling Technique (SMOTE) [6] to solve data imbalance in EDMS datasets. The main innovation is an ensemble-based model that improves the accuracy of EDMS identification by combining Principal Component Analysis (PCA) with AdaBoost, with a focus on the critical 2-4 months postpartum period. In the end, impacted children and their families will benefit from this model's superior performance over traditional approaches, which offers the possibility of better early detection and intervention in EDMS while maintaining academic integrity.

**1.1.3 Scope and significance of the study**

Encephalitic Disintegration Morbidity Speculations (EDMS) is a neurodevelopmental condition that usually appears in early childhood. The main goal of this work is to improve the accuracy of EDMS identification by utilizing sophisticated machine-learning techniques. Through a thorough comparative analysis of multiple machine learning algorithms, the research aims to determine which algorithm produces the best results consistently. Additionally, the Synthetic Minority Over-Sampling Technique (SMOTE) is employed to address the ongoing problem of data imbalance in EDMS datasets. One noteworthy development is the creation of an ensemble-based predictive model that makes use of Principal Component Analysis (PCA)[7] and AdaBoost to greatly increase the accuracy of EDMS detection, especially in the critical window of 2-4 months after childbirth, allowing for early identification and intervention. In addition to expanding our knowledge of EDMS and machine learning, this research has the potential to transform early EDMS diagnosis and care, minimizing the long-term effects on patients and their families.

**LITERATURE SURVEY**

The research conducted by Taban Eslami, Vahid Mirjalili, Alvis Fong, Angela R. Laird, and Fahad Saeed [3] presents "ASD-DiagNet," an innovative approach for the accurate diagnosis of autism spectrum disorder (ASD), with a particular focus on pediatric patients. This method exclusively utilizes fMRI data and employs a sophisticated joint learning technique to optimize feature extraction. Furthermore, it incorporates a distinctive data augmentation strategy, leading to remarkable improvements in classification accuracy and efficiency, streamlining analysis processes, and significantly reducing execution time.

Thabtah et al. [1] used Information Gain (IG) and Chi-Squared (CHI) approaches to develop feature subsets specifically designed for adults and adolescents in a different study. Logistic regression (LR) was then used to integrate these subgroups into the diagnosis procedure for autism spectrum disorder (ASD). By introducing forward feature selection and under-sampling methods, the study improved the accuracy of ASD diagnosis for people and advanced the development of diagnostic methodologies related to ASD. This study is a reflection of the continuous efforts to improve ASD diagnosis by utilizing cutting-edge data analysis and feature engineering methods.

The paper by Firuz Kamalov and Fadi Thabtah[2] discusses Autism Spectrum Disorder (ASD) and the requirement for effective screening instruments. Diagnoses with ASD are rising, which frequently causes examinations to be postponed. They present Variable Analysis (Va), a computational intelligence technique that minimizes feature correlations while identifying important aspects of ASD screening tools. Using machine learning techniques, they assess Va's performance using various criteria, including predicted accuracy, sensitivity, and specificity. The findings are encouraging: Va preserves predictive accuracy while lowering the quantity of screening items required for adults, teenagers, and toddlers. This strategy emphasizes the use of intelligent techniques and machine learning to boost ASD screening precision and efficiency while providing simplified, user-friendly ASD screening instruments for early identification and better results.

Sankar K. Pal and Sushmita Mitra [4] present a novel neural network model for fuzzy pattern classification in their study, "Multilayer Perceptron, Fuzzy Sets, and Classification." This novel model tackles the problems posed by imprecise or ambiguous data. It includes uncertainty into the classification process efficiently by using linguistic input representation and basing output judgments on class membership values. The study highlights the potential advantages of merging neural network-based categorization with fuzzy notions and indicates that process efficiency might be greatly increased by implementing parallel hardware.

Reem Ahmed Bharathiq [5] and colleagues discuss the difficulties in diagnosing autism spectrum disorder (ASD), a complicated neurodevelopmental disorder that affects about 1% of the population, in their work on autism spectrum disorder (ASD) diagnosis using structural MRI and machine learning. By examining structural MRI data, they investigate how machine learning might enhance diagnosis. Their research shows the potential of machine learning (ML) in improving our knowledge of autism spectrum disorders (ASD) and creating tailored diagnostic tools for medical professionals, despite some obstacles such as limited sample sizes.

**PROPOSED SOLUTION**

The suggested method makes use of machine learning to identify autism in young children, including adults and adolescents as well as low-functioning toddlers with Encephalitic Disintegration Morbidity Speculations (EDMS) between the ages of two and four months. AdaBoost is an efficient boosting technique that uses iterative learning to decrease misclassifications and increase accuracy. The study looks at a variety of EDMS features and compares several algorithms; AdaBoost stands out for correctly diagnosing people with EDMS traits.
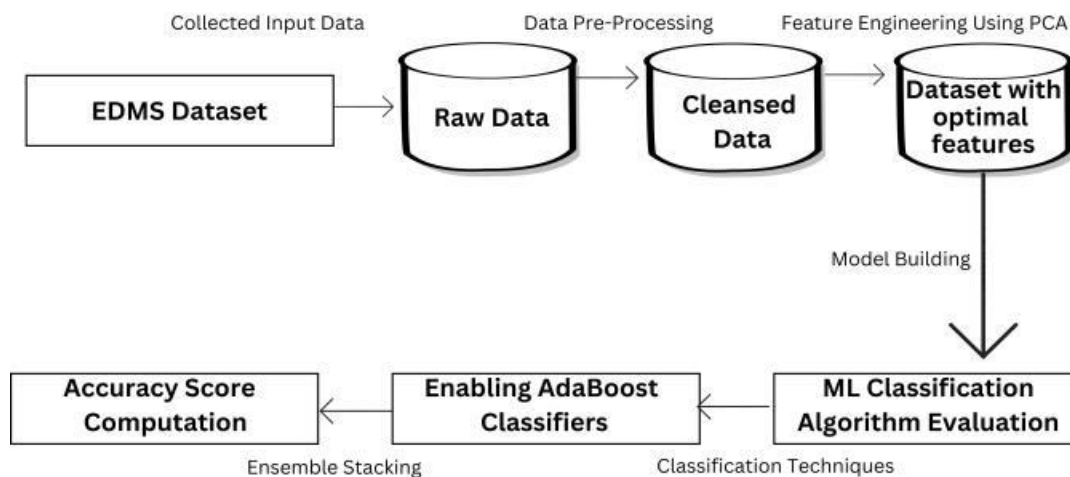
**Figure 1:** Architecture diagram of this Proposed solution to predict the accuracy of the disorder

### 3.1.1 Data Collection and Pre-processing

An extensive dataset covering essential characteristics is gathered to discover Encephalitic Disintegration Morbidity Speculations (EDMS) in infants between the ages of two and four months early. This dataset assesses a child's growth and health by looking at milestones in development, physical activity levels, upper limb movement, and childhood obesity. Data preparation is the next step, which addresses outliers, missing numbers, and inconsistencies to guarantee data correctness and consistency. To ensure uniform feature scaling, standardization and normalization procedures are used, which offer a strong basis for precise EDMS identification.

The study makes use of the Encephalitic Disintegration Morbidity Speculations (EDMS) Screening Datasets from the UCI database, which are divided into various age groups. There are many different kinds of data in these datasets, such as binary, continuous, and category data. First, attributes containing "Null Value," "Irrelevant Values," and "Qchat-10-Score" are eliminated. This produces a refined dataset that is appropriate for modeling and in-depth analysis.

Table 3.1: Datasets Summary

| Dataset | Sample Size | Feature Size with Class value | Classes | Presence of Missing Attributes |
|---|---|---|---|---|
| ASD Screening Data For Children | 292 | 21 | ASD/Non-ASD | Yes |
| Autism & Neuro developmental Service For Toddlers | 704 | 21 | ASD/Non-ASD | yes |
| ASD Screening For Adolescent | 104 | 21 | ASD/Non-ASD | yes |
| ASD Screening For Adult | 1054 | 21 | ASD/Non-ASD | yes |

### 3.1.2. Feature Selection

Isolating and incorporating relevant attributes through effective feature selection is essential for improving the performance of machine learning models. Sophisticated feature selection techniques have been used with consideration to a variety of datasets featuring adults, adolescents, toddlers, and children in the context of EDMS identification. Recursive Feature Elimination (RFE), Correlation-based Feature Selection with Harmony Search (CFS–Harmony Search), and the Boruta algorithm are three of the well-known methods that stand out among the others.

**3.1.2.1 Boruta Algorithm**

The wrapper approach that is being provided extends the random forest (RF) technique, which is meant for feature selection. To generate shadow features, each feature in the dataset is copied, and their values are then randomly combined. Subsequently, the method employs RF on the revised dataset to ascertain the importance of every characteristic. The highest real feature score is compared to the highest shadow feature score for every iteration once both real and shadow feature scores have been calculated. During this process, features deemed to be of minimal importance are carefully eliminated. When all features are accepted or rejected, or when a certain number of RF runs is reached, the algorithm comes to an end.

**3.1.2.2 Correlation-Based Feature Selection with Harmony Search**

Relationship-Based A heuristic function is used in feature selection using harmony search to order features according to their correlations. It takes into account variables such as the number of indicators (M), the number of harmonies (N) that are stored in memory, and their potential values (D). It also takes into account the ideal indication (i) in the harmony memory (Ei) and the harmony memory rate (Et). The formula

$$P = Ei / D \times (1 - Et) \qquad (3.1)$$

is used to calculate the probability (P). This technique leverages feature correlations to find an ideal subset and optimize different parameters for better outcomes. It combines feature selection with harmony search optimization.

**3.1.2.3 Recursive Feature Elimination (RFE)**

Recursive Feature Elimination (RFE) is a methodical feature selection strategy wherein the least significant characteristics are eliminated progressively. The process starts with the classifier being trained using all of the initial characteristics. Each feature's significance is evaluated according to how well it fits the job at hand. The feature with the lowest score, which denotes a lower level of relevance, is then removed from the original feature set. Until the required number of feature subsets is reached, this recursive process is repeated recursively. The RFE process can be summed up as follows: 1. Training of the Classifier: The entire set of features is initially used to train the classifier. 2. The process of calculating a score involves quantifying each feature's relevance and rating them. 3. Feature Elimination: A feature's reduced relevance is indicated by its lowest score, which is systematically removed from consideration.

**3.1.3 Applying Individual Classification Algorithms**

The current study comprised the methodical creation of feature subsets, which were subsequently tested using several classifiers. Classifiers that did not meet the expected performance were excluded. Using the remaining classifiers—Random Tree (RT), K-Nearest Neighbours (KNN), Support Vector Machine (SVM), Naive Bayes (NB), Multi-layer Perceptron (MLP), Logistic Regression (LR), and Ada Boost—a comparative analysis was conducted. The best-performing classifier was identified and the feature subsets associated with its peak performance were identified based on specific metrics. Using this methodology simplified the process of selecting the most informative feature subsets and the best classifier in an orderly fashion for the given problem.

**3.1.3.1 Naïve Bayes (NB)**

Based on Bayes' Naïve Bayes (NB) is an incredibly powerful probabilistic classifier.

$$R(b/a) = R(a/b)/R(b) \qquad (3.2)$$

R(a | b) in the context of NB denotes the posterior probability of the target class, whereas R(a) denotes the probability of the prior class. R(b) represents the predictor's prior probability, while R(b | a) indicates the likelihood of the supplied predictor about a specific class. By estimating the likelihood of a particular class based on the combined probabilities of several predictor factors, NB's fundamental idea is to efficiently predict target outputs. The simplicity of this classifier and its reliance on the conditional independence of the predictor variables make it especially useful for producing accurate and timely classification results.

**3.1.3.2 Support Vector Machine (SVM)**

Data points in a multi-dimensional space can be divided into distinct classes using the Support Vector Machine (SVM), a potent algorithm that creates a decision boundary known as a hyperplane. To reduce classification errors and increase the distance between classes, SVM typically builds two parallel hyperplanes. By finding the ideal balance between accuracy and the separation between data points of different classes, this configuration often referred to as a maximum margin classifier ensures the most reliable and ideal categorization.

**3.1.3.3 K- Nearest Neighbour (KNN)**

K-Nearest Neighbour (KNN) is a method used for example categorization that primarily depends on how close data points are to each other. It often takes into account several neighbours and uses the Euclidean approach, which is represented by the following equation, to calculate their distances:

$$D = \sqrt{(m_1 - m_2)^2 + (n_1 - n_2)^2} \qquad (3.3)$$

The distance between two points, indicated as (m1, n1) and (m2, n2), is represented by the letter "D" in this equation. By evaluating the classes of a data point's closest neighbors, KNN is very helpful in determining the class of the data point. Pattern recognition, classification, and data analysis can all benefit from this proximity-based method.

**3.1.3.4 Random Tree (RT)**

A decision tree variant known as Random Tree (RT) creates several alternative trees in a stochastic way, each with a subset of K random attributes. By using ensemble learning, this method improves prediction accuracy and generates reliable forecasts. It generates different tree designs by adding randomness to feature selection, combining their results to enhance prediction

performance overall. In machine learning, random trees are frequently employed to handle complicated, high-dimensional data in a variety of classification and regression tasks.

### 3.1.3.5 Logistic Regression (LR)

Logistic Regression (LR) is a statistical technique designed for handling binary dependent variables, making it particularly suitable for binary classifications where output values fall within the 0 to 1 range, such as true or false and yes or no. Unlike standard linear regression, LR can effectively manage datasets that contain nominal, categorical, or discrete variables. It establishes the relationship between one or more nominal or ordinal predictor variables and a single binary dependent variable by employing a sigmoidal function to transform input variables into a probability score. Logistic Regression plays a vital role in various fields like economics, epidemiology, and machine learning, where accurately estimating the probability of binary outcomes is essential for decision-making and modeling.

### 3.1.3.6 AdaBoost (Adaptive Boosting)

An ensemble classifier called AdaBoost (Adaptive Boosting) lowers errors by repeatedly combining several weak classifiers. Based on their accuracy, it gives training samples and classifiers weights. It trains weak classifiers on random selections of data throughout each cycle, and then merges the results. The system improves overall classification performance by prioritizing difficult situations and dynamically adjusting weights.

## RESULT AND DISCUSSION

### 4.1.1 Data Set

Real datasets from Kaggle and the UCI repository were used to thoroughly assess the EDMS model, with an emphasis on a range of age groups, including infants, toddlers, adolescents, and adults. Each dataset included a large number of samples, with 292,704, 104, and 1,054 cases for children, toddlers, adolescents, and adults, respectively, and 21 unique attributes. Table [1] provides a thorough collection of variable details and dataset acronyms. The dependability and efficiency of the EDMS model are confirmed by this thorough evaluation, which is based on a wide range of age groups and sizable sample sizes. As a result, the EDMS model is a valuable addition to the area of study and a good contender for publishing in scholarly publications.

### 4.1.2 Performance Metrics

Utilizing the below measures listed, the machine learning model's performance is measured. For this evaluation, the number of correct forecasts and the percentage of forecasts across various threshold values are examined. Understanding how well the model captures different prediction scenarios and maximizes its performance is made easier with this technique. Key performance parameters, including Precision, Recall, F1-score, and Accuracy, are used in the model review process.

### 4.1.2.1 Accuracy

The percentage of accurately anticipated data points relative to the total number of predictions is known as accuracy, and it is a basic statistic. It is noteworthy that the EDMS model achieves about 100% accuracy in diagnosing the illness.
The following is the formula used to determine accuracy:
Accuracy = No of Predictions / Total No of Predictions

### 4.1.2.2 Precision

Precision can be defined as the ratio of the total number of actual positive predictions to the number of correct positive predictions made by the EDMS proposed model. The model achieved a flawless precision value of 1.0, indicating that it obtained perfection. Precision is computed using the following formula:
Precision = No of Correct Positive Prediction \ Total No of Positive Prediction

### 4.1.2.3 Recall

The ratio of all positive cases to the model's correctly predicted positive outcomes is known as recall. Another excellent and noteworthy result is the recall value of 1.0 that was attained. This is the recall formula:
Recall=No of Correct Positive Prediction\Total No of Values in Positive Predictions

### 4.1.2.4 F1-Score

The precision and recall values are weighted to create the F1-score, and the resultant F1-score is 1.0. The F1-score calculation formula is:
F1-Score is equal to 2 * (Recall * Precision) / (Recall + Precision).
The model performs exceptionally well, as evidenced by these remarkable precision, recall, and F1-score values, making it a solid and trustworthy diagnostic tool for well-defined identification of disorder.
By contrasting different single-learning classifiers, a new model with improved early autism prediction potential is shown in Figure [2]. Accuracy and execution time can be enhanced by utilizing Principal Component Analysis (PCA) and integrating the AdaBoost Classifier into the Ensemble Classification for Autism pre-diagnosis (EDMS). The outcomes of the EDMS model surpassed those of earlier models, providing a useful instrument for the early diagnosis of autism spectrum disorder (ASD).
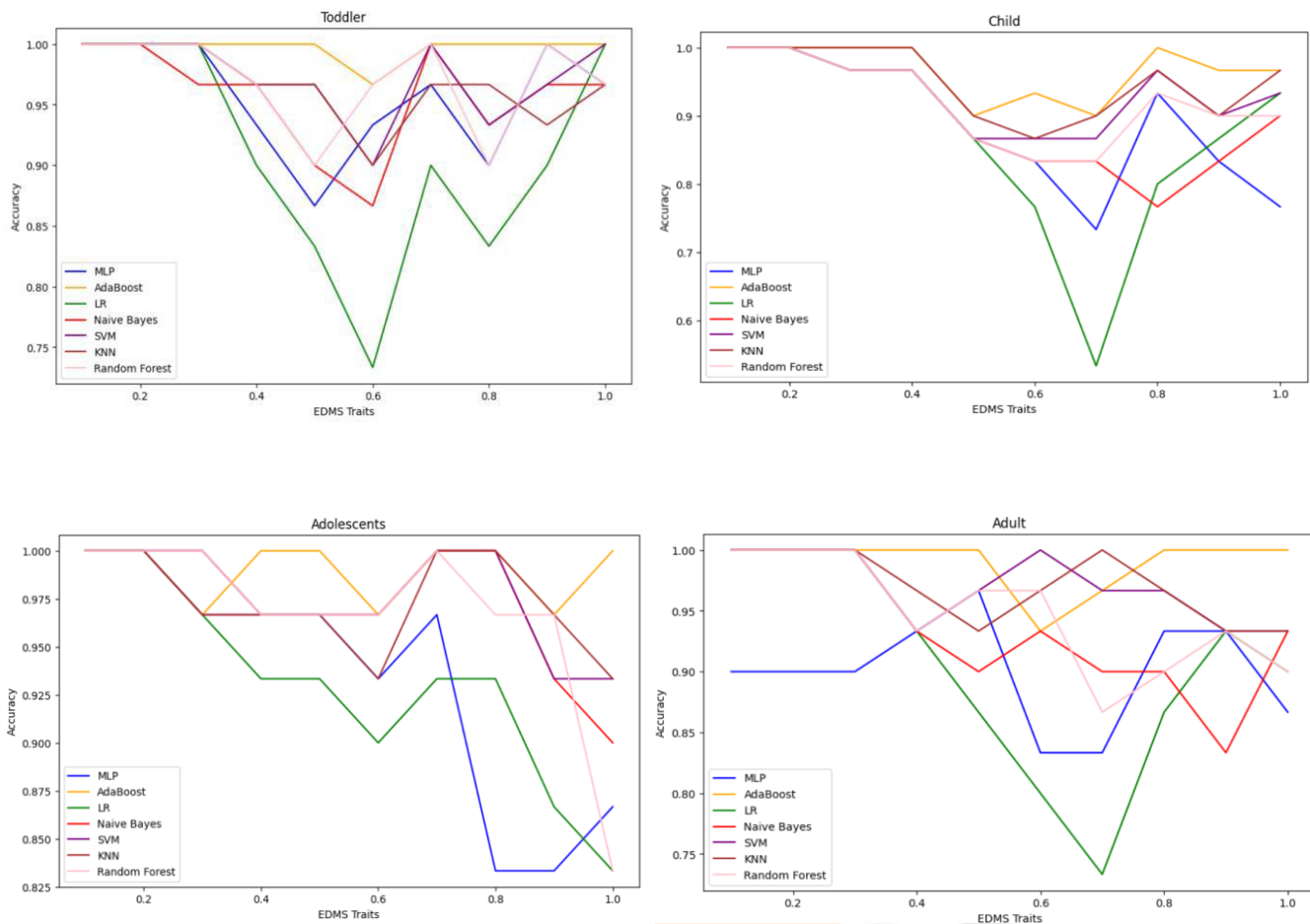
Figure 2. Accuracy Comparison Among Different ML Models with Toddler, Child, Adolescents, and Adult Datasets.

## CONCLUSION

In summary, EDMS can appear early in life and have a significant impact on a person's physical, social, emotional, and cognitive health. The difficulties in correctly diagnosing and treating EDMS highlight how crucial this research is to improving our comprehension of the complex aspects of this illness. Principal Component Analysis (PCA) was utilized to predict EDMS, and it was contrasted with several different techniques, such as Support Vector Machines (SVM), Naive Bayes, Random Forest (RF), Logistic Regression (LR), KNN, and AdaBoost. It was evident from this comparison that AdaBoost yielded the most accurate findings. Notably, the Existing method also investigated the application of unsupervised machine learning methods to autism prediction. These techniques did not, however, produce the required level of accuracy. Accuracy was greatly increased by incorporating the Synthetic Minority Over-Sampling Technique (SMOTE) to overcome this. The target demographic expansion is one significant improvement in the suggested methodology. The method goes beyond the current system's limitations by forecasting EDMS for children, adults, and adolescents in addition to adults, toddlers, and adolescents. With this extension, people will have the chance for early intervention and diagnosis, which is a big development in early-stage diagnosis. The future direction of this research will center on using larger datasets along with early detection in the fetal stage to improve accuracy even more. Even more accurate predictions and a deeper comprehension of EDMS are anticipated as a result of access to larger datasets, which the study is currently using to work with.

## ACKNOWLEDGMENT

## REFERENCES

[1] Thabtah F, Kamalov F, Rajab K. A new computational intelligence approach to detect autistic features for autism screening. Int J Med Inform. 2018 Sep; 117:112-124. doi: 10.1016/j.ijmedinf.2018.06.009. Epub 2018 Jun 27. PMID: 30032959.
[2] Eslami T, Mirjalili V, Fong A, Laird AR, Saeed F. ASD-DiagNet: A Hybrid Learning Approach for Detection of Autism Spectrum Disorder Using fMRI Data. Front Neuroinform. 2019 Nov 27; 13:70. doi: 10.3389/fninf.2019.00070. PMID: 31827430; PMCID: PMC6890833.

[3] Moridian P, Ghassemi N, Jafari M, Salloum-Asfar S, Sadeghi D, Khodatars M, Shoeibi A, Khosravi A, Ling SH, Subasi A, Alizadehsani R, Gorriz JM, Abdulla SA, Acharya UR. Automatic autism spectrum disorder detection using artificial intelligence methods with MRI neuroimaging: A review. Front Mol Neurosci. 2022 Oct 4; 15:999605. doi: 10.3389/fnmol.2022.999605. PMID: 36267703; PMCID: PMC9577321.

[4] Hossain MD, Kabir MA, Anwar A, Islam MZ. Detecting autism spectrum disorder using machine learning techniques: An experimental analysis on toddler, child, adolescent and adult datasets. Health Inf Sci Syst. 2021 Apr 6;9(1):17. doi: 10.1007/s13755-021-00145-9. PMID: 33898020; PMCID: PMC8024224.

[5] Bahathiq RA, Banjar H, Bamaga AK, Jarraya SK. Machine learning for autism spectrum disorder diagnosis using structural magnetic resonance imaging: Promising but challenging. Front Neuroinform. 2022 Sep 28; 16:949926. doi: 10.3389/fninf.2022.949926. PMID: 36246393; PMCID: PMC9554556.

[6] Ismail E, Gad W, Hashem M. A hybrid Stacking-SMOTE model for optimizing the prediction of autistic genes. BMC Bioinformatics. 2023 Oct 6;24(1):379. doi: 10.1186/s12859-023-05501-y. PMID: 37803253; PMCID: PMC10559615.

[7] Shihab AI, Dawood FA, Kashmar AH. Data Analysis and Classification of Autism Spectrum Disorder Using Principal Component Analysis. Adv Bioinformatics. 2020 Jan 7; 2020:3407907. doi: 10.1155/2020/3407907. PMID: 32395129; PMCID: PMC7199592.

[8] O. Altay and M. Ulas, "Prediction of the autism spectrum disorder diagnosis with linear discriminant analysis classifier and K-nearest neighbor in children," 2018 6th International Symposium on Digital Forensic and Security (ISDFS), Antalya, Turkey, 2018, pp. 1-4, doi: 10.1109/ISDFS.2018.8355354.

[9] Alves CL, Toutain TGLO, de Carvalho Aguiar P, Pineda AM, Roster K, Thielemann C, Porto JAM, Rodrigues FA. Diagnosis of autism spectrum disorder based on functional brain networks and machine learning. Sci Rep. 2023 May 18;13(1):8072. doi: 10.1038/s41598-023-34650-6. PMID: 37202411; PMCID: PMC10195805.

[10] Bi XA, Wang Y, Shu Q, Sun Q, Xu Q. Classification of Autism Spectrum Disorder Using Random Support Vector Machine Cluster. Front Genet. 2018 Feb 6; 9:18. doi: 10.3389/fgene.2018.00018. PMID: 29467790; PMCID: PMC5808191.

[11] Demirhan, A., "Performance of machine learning methods in determining the autism spectrum disorder cases", Mugla Journal of Science and Technology, 4(1), 79-84, 2018.

[12] Brugha TS, McManus S, Smith J, et al. Validating two survey methods for identifying cases of autism spectrum disorder among adults in the community. Psychol Med 2012; 42:647–56.