

A FIELD PROJECT REPORT

on

**“Unlocking Autism DNA”**

**Submitted**

by

221FA04140  
K Jhansi Suvarchala

221FA04113  
S Mercy

221FA04143  
Sk Mahmooda Aafreen

221FA04074  
Ch Sowmya

**Under the guidance of**

*Maridu Bhargavi*

*Assistant Professoress, Department of CSE, VFSTR*



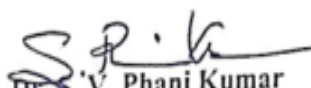
**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**  
**VIGNAN'S FOUNDATION FOR SCIENCE, TECHNOLOGY AND RESEARCH**  
**Vadlamudi, Guntur.**  
**ANDHRA PRADESH, INDIA, PIN-522213.**

### **CERTIFICATE**

This is to certify that the Field Project entitled “**UNLOCKING AUTISM DNA**” that is being submitted by 221FA04140 (K Jhansi Suvarchala), 221FA04113(S Mercy), 221FA04143(Sk Mahmooda Aafreen), 221FA04074 (CH Sowmya) for partial fulfilment of Field Project is a bonafide work carried out under the supervision of M Bhargavi , M.Tech., Associate Professor, Department of CSE.

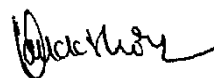
M Bhargavi

Assistant Professoress, CSE



Dr. S. V. Phani Kumar

HOD,CSE



Dr.K.V. Krishna Kishore

Dean, SoCI



## DECLARATION

We hereby declare that the Field Project entitled “**UNLOCKING AUTISM DNA**” that is being submitted by 221FA04140 (K Jhansi Suvarchala), 221FA04113(S Mercy), 221FA04143(Sk Mahmooda Aafreen), 221FA04074 (CH Sowmya) in partial fulfilment of Field Project course work. This is our original work, and this project has not formed the basis for the award of any degree. We have worked under the supervision M Bhargavi , Associate Professor, Department of CSE

By

221FA04140 (K Jhansi Suvarchala)

221FA04113(S Mercy)

221FA04143(Sk Mahmooda Aafreen)

221FA04074 (Ch Sowmya)

Date:

# ABSTRACT

Autism Spectrum Disorder (ASD) is a neurodevelopmental condition diagnosed through various behavioral assessments. Early diagnosis is crucial for effective intervention and care, but traditional clinical methods can be time-consuming and prone to errors. To address these limitations, machine learning models like Logistic Regression, Support Vector Machines (SVM), XGBoost, and AdaBoost have been applied for earlier detection of ASD. These models can offer high accuracy in predicting ASD, but there is a significant concern regarding the security and privacy of sensitive patient data in previous research. Our study evaluates the performance of machine learning models, such as Logistic Regression, AdaBoost, and ensemble techniques like XGBoost and Random Forest. While these models have achieved excellent accuracy in some cases, a major drawback in existing studies is their vulnerability to data breaches, which compromises patient confidentiality. To address this issue, we propose using Federated Learning (FL), which offers a more secure and privacy-preserving solution. FL enables models to be trained collaboratively across decentralized devices without directly sharing patient data, thus mitigating privacy risks. This decentralized approach ensures local processing of sensitive information, allowing for enhanced data security while still improving the overall performance of the models through collaborative learning. By integrating FL, we aim to provide a more secure and efficient solution for the early detection of ASD, balancing both accuracy and privacy in machine learning applications for healthcare.

## TABLE OF CONTENTS

| Section | Title                                       | Page no |
|---------|---|---------|
| 1.      | Introduction                                | 1-3     |
| 1.1     | What is Autism Spectrum Disorder (ASD)?     | 2       |
| 1.2     | The Importance of Early Diagnosis           | 2       |
| 1.3     | Current Machine Learning Techniques for ASD | 2       |
| 1.4     | Limitations of Traditional Models           | 3       |
| 1.5     | Federated Learning as a Solution            | 3       |
| 2.      | Literature Survey                           | 4-7     |
| 2.1     | Literature Review for ASD Prediction        | 4-7     |
| 2.2     | Motivation                                  | 7       |
| 3.      | Methodology                                 | 8-14    |
| 3.1     | Statement of the Problem                    | 9       |
| 3.2     | Data Collection                             | 9       |
| 3.3     | Exploratory Data Analysis                   | 9-11    |
| 3.4     | Data Preprocessing                          | 11      |
| 3.4.1   | Missing Values                              | 11      |
| 3.4.2   | Normalization                               | 11      |
| 3.5     | Federated Learning Approach                 | 11      |
| 3.5.1   | Node Training                               | 11      |
| 3.5.2   | Model Update                                | 11      |
| 3.5.3   | Application of PCA                          | 12      |
| 3.6     | Train-Test Split                            | 12      |
| 3.7     | Application of Machine Learning Algorithms  | 12- 14  |
| 3.7.1   | Logistic Regression                         | 12      |
| 3.7.2   | K-Nearest Neighbors (KNN)                   | 12      |

| <b>Section</b> | <b>Title</b>                                    | <b>Page no</b> |
|----------------|---|----------------|
| 3.7.3          | Support Vector Machines (SVM)                   | 13             |
| 3.7.4          | Decision Trees                                  | 13             |
| 3.7.5          | Naive Bayes                                     | 13             |
| 3.7.6          | AdaBoost  | 14             |
| 3.7.7          | Ensemble (XGBoost + Random Forest)              | 14             |
| 3.8            | Model Evaluation                                | 14             |
| 4.             | Implementation                                  | 15 -1 7        |
| 4.1            | Environment Setup                               | 16             |
| 4.2            | Sample Code for Preprocessing and ML Operations | 16 – 17        |
| 5.             | Results and Analysis                            | 18- 25         |
| 5.1            | Performance Metrics of ML Models                | 19             |
| 5.2            | Confusion Matrices                              | 19-23          |
| 5.3            | Comparison with Previous Models                 | 25             |
| 6.             | Conclusion                                      | 26-27          |
| 7.             | References                                      | 28-31          |

## LIST OF FIGURES

|   |    |
|---|----|
| <b>Figure 1:</b> Flowchart of Methodology                               | 9  |
| <b>Figure 2:</b> Pie chart showing 69.1% of people with ASD             | 10 |
| <b>Figure 3:</b> ASD cases by sex, showing a higher prevalence in males | 11 |
| <b>Figure 4:</b> ASD cases and jaundice prevalence                      | 11 |
| <b>Figure 5:</b> ASD cases by family members with ASD                   | 11 |
| <b>Figure 6:</b> Comparison of model accuracies                         | 19 |
| <b>Figure 7:</b> Confusion Matrix - Logistic Regression                 | 21 |
| <b>Figure 8:</b> Confusion Matrix – KNN                                 | 21 |
| <b>Figure 9:</b> Confusion Matrix – SVM                                 | 22 |
| <b>Figure 10:</b> Confusion Matrix - Decision Tree                      | 22 |
| <b>Figure 11:</b> Confusion Matrix - Naive Bayes                        | 23 |
| <b>Figure 12:</b> Confusion Matrix – Adaboost                           | 23 |
| <b>Figure 13:</b> Confusion Matrix - Ensemble Method (XGB + RF)         | 24 |

## LIST OF TABLES

|  |    |
|--|----|
| <b>Table I:</b> Model Accuracy Comparison                | 20 |
| <b>Table II:</b> Model Performance Metrics in Percentage | 25 |



# **CHAPTER 1**

## **INTRODUCTION**

# **1. INTRODUCTION**

## **1.1 What is Autism Spectrum Disorder (ASD)?**

Autism Spectrum Disorder (ASD) is a neurodevelopmental condition characterized by challenges in communication, social interactions, and repetitive behaviors. ASD manifests differently in each individual, making diagnosis difficult. It is considered multifactorial, involving both genetic and environmental factors, which contribute to its wide range of symptoms. Traditional diagnosis methods are often subjective and time-consuming, relying heavily on behavioral observations. These complexities make early detection challenging, but essential, as timely intervention can significantly improve outcomes for individuals affected by ASD.

## **1.2 The Importance of Early Diagnosis**

Early diagnosis of Autism Spectrum Disorder (ASD) is critical as it allows for timely interventions that significantly improve the individual's quality of life. Early therapeutic support helps address developmental delays and social challenges faced by those with ASD. Unfortunately, traditional diagnostic methods often delay early detection due to reliance on behavioral observations, which are typically recognized only after two years of age. Machine learning models offer an opportunity for earlier, more objective diagnosis by analyzing behavioral, genetic, and environmental data, leading to earlier intervention and better developmental outcomes for children with ASD.

## **1.3 Current Machine Learning Techniques for ASD**

Machine learning (ML) techniques, such as Logistic Regression, SVM, XGBoost, and AdaBoost, have proven effective in identifying Autism Spectrum Disorder (ASD) by analyzing complex datasets. These models can predict ASD with high accuracy, offering a significant improvement over traditional diagnostic methods, which rely on subjective behavioral assessments. ML enables quicker, more objective analysis of ASD risk factors, aiding in early diagnosis. Despite their accuracy, these models face challenges with data privacy and security, as they often require access to sensitive health information, which makes their deployment in healthcare settings more difficult.

#### **1.4 Limitations of Traditional Models**

Traditional machine learning models for Autism Spectrum Disorder (ASD) detection face several limitations, primarily related to privacy and security concerns. Most models rely on centralized datasets that require sensitive patient health information, making them vulnerable to data breaches. Additionally, these models often suffer from biases in training data, limiting their generalizability across diverse populations. Inconsistent preprocessing methods and smaller datasets further restrict their effectiveness. These challenges make it difficult for traditional models to provide reliable and secure early ASD detection, highlighting the need for more secure and privacy-conscious approaches.

#### **1.5 Federated Learning as a Solution**

Federated Learning (FL) addresses privacy concerns in traditional machine learning by allowing data to remain on local devices, only sharing model updates like gradients with a central server. This ensures that sensitive patient data is never exposed, protecting privacy while still enabling collaborative learning across multiple devices. For Autism Spectrum Disorder (ASD) diagnosis, FL enables secure, decentralized data processing, making it an ideal solution for maintaining confidentiality. FL overcomes the limitations of centralized models, ensuring that ASD detection remains accurate, while safeguarding patient privacy through its decentralized, secure data-handling approach.

# **CHAPTER-2**

## **LITERATURE SURVEY**

## 2. LITERATURE SURVEY

### 2.1 Literature review

We have carried out a literature survey to include all related works with our study on Autism Prediction. The crux of ideas from these papers has been summed up below:

| No | Author(s)             | Model/Approach  | Accuracy/Results                   | Limitation  |
|----|-----------------------|---|------------------------------------|---|
| 1  | N. BalaKrishna et al. | SVM, Decision Trees, Linear Discriminant Analysis, and Logistic Regression                                    | <b>SVM:</b> 93%                    | Effective pre-processing methods and high- quality data are needed to boost the SVM model's performance.  |
| 2  | N. Zaman et al.       | Logistic Regression, KNN, SVC, Naive Bayes, Decision Tree and Random Forest                                   | <b>Naive Bayes:</b> 96.23%         | not cover all autism spectrum variations and might require validation across diverse populations to ensure broad applicability.                                     |
| 3  | S. Islam et al.       | K-NN, Decision Tree, Random Forest, SVM, Logistic Regression, Naive Bayes and Gradient Boosting               | <b>KNN:</b> 98%                    | The only limitation of our model is the lack of enough large data to train our model.   |
| 4  | S. K. R. Naik et al.  | KNN, Logistic Regression, Decision Tree Classifier, Random Forest Classifier, Naive Bayes, XGBoost Classifier | <b>XGBoost Classifier:</b> 98.2%   | it may not address potential limitations such as dataset bias, generalizability, or the applicability of the system to diverse populations or real-world scenarios. |
| 5  | Vakadkar et al.       | SVM, RFC, NB, LR and KNN  | <b>Logistic Regression:</b> 97.15% | The primary limitation of this research is the scarce availability of large and open source ASD datasets. Does not predict severity of ASD                          |
| 6  | A. Baranwal et al.    | LDA, NB, Classification And Regression Trees (CART), KNN, LR and SVM  | <b>LDA:</b> 72.2024%               | The datasets especially the child and adolescent datasets, are really small in size and are not suitable for building machine learning models.                      |

|    |                                |   |   |   |
|----|--------------------------------|---|---|---|
| 7  | R. Chauhan et al.              | SVM, KNN, Random Forest, Decision Tree  | <b>Random Forest:</b> 74%   | It needs a wide variety of computing infrastructures, such as real-time tracking, collaborative tools, and statistical analysis.  |
| 8  | Aishwarya D et al.             | Decision Tree, Random Forest, KNN, SVM, Logistic Regression, Naive Bayes, Gradient Boosting, K-Means Clustering, PCA and Neural Networks. | <b>Neural Networks:</b> 90-99%  | results in less expensive therapy and better patient outcomes. The development of the field of autism research will benefit from the use of this dashboard.                         |
| 9  | Y. J. Cheong et al.            | Logistic Regression, SVM, Decision Trees, Random Forests, KNN, and Gradient Boosting Machines.  | <b>Random Forests</b> - 97%   | the study does not address the limitations associated with the model's applicability across diverse ASD populations or the generalizability of the findings beyond the sample used. |
| 10 | K. -F. Kollias et al.          | the application of robotics and artificial intelligence in ASD diagnosis and emotion recognition.   | robotics has significant potential in enhancing ASD diagnosis and emotional support.  | Some studies had small sample sizes, and the effectiveness of robots varied, indicating a need for more rigorous research and broader technology integration.                       |
| 11 | S. Bose et al.                 | Logistic Regression, SVC, XGBoost and Naive Bayes   | XGBoost up to 100% accuracy   | their effectiveness is contingent on the quality of the datasets used, which can affect sensitivity and specificity.  |
| 12 | C. K. Syriopoulou-Delli et al. | RNN, SVM, DNN and Linear classifiers  | Long Short-Term Memory (LSTM) networks and SVM, achieves high classification accuracy | the variability in sample sizes, age groups, and functional skills of participants, which can affect the generalizability of findings.  |
| 13 | S. Kumar et al.                | CNNs on MRI brain scans   | 93-94%  | Requires curated MRI datasets and high computational cost.  |
| 14 | Y. Wang et al.                 | VGG16, ResNet for Facial Recognition  | 90%   | Privacy concerns and false positives.   |
| 15 | A. Sharma et al.               | Federated Learning with Smart Devices   | 95%   | Limited by device interoperability and secure data transmission.  |
| 16 | P. Gupta et al.                | ASD-EVNet: Ensemble Vision Network  | 95%   | False positives in non-ASD behavior.  |
| 17 | M. Chen et al.                 | Self-Attention DNN on fMRI Data   | 96%   | Expensive and requires brain imaging datasets.  |

|    |                    |  |      |   |
|----|--------------------|--|------|---|
| 18 | R. Patel et al.    | Eye-Tracking and Scanpath Analysis             | 92%  | Needs specialized hardware.   |
| 19 | F. Ali et al.      | NLP-Based Speech Pattern Recognition           | >90% | Drops in performance with noisy or missing speech data.               |
| 20 | T. Bose et al.     | Logistic Regression on M-CHAT Data             | 85%  | Simple model unsuitable for complex ASD cases.                        |
| 21 | L. Zhang et al.    | Random Forest with Genetic Algorithm           | 95%  | Computationally intensive feature selection.                          |
| 22 | K. Naik et al.     | Hybrid XGBoost and SVM Models                  | 96%  | Requires careful tuning.  |
| 23 | S. Dutta et al.    | LSTM on Behavioral Time-Series Data            | 94%  | Requires long-term behavioral data collection.                        |
| 24 | A. Baranwal et al. | Autoencoders for Anomaly Detection             | 90%  | Depends heavily on data quality.                                      |
| 25 | C. Lee et al.      | Federated Learning for Collaborative Diagnosis | 97%  | Privacy concerns and challenges in collaboration across institutions. |

## 2.2 Motivation

Autism cannot define a person; however, the personality of each with autism is unique. Acceptance and understanding may be the way to care for people with autism. Individuals with autism can have major positive outcomes if identified early; Autism isn't illness; moreover, nobody requires a cure to possess autism; people with autism have various strengths and abilities which may be capitalized on to produce a lot of things. Autism is not a limitation but the way of thinking and experiencing the world. Individuals with autism will then, with proper support and accommodations, bloom and flourish in all areas of life. Detection of autism is no labeling or stigmatizing but providing support to succeed. Every individual diagnosed with autism has the potential to be a good influence in this world.

# **CHAPTER-3**

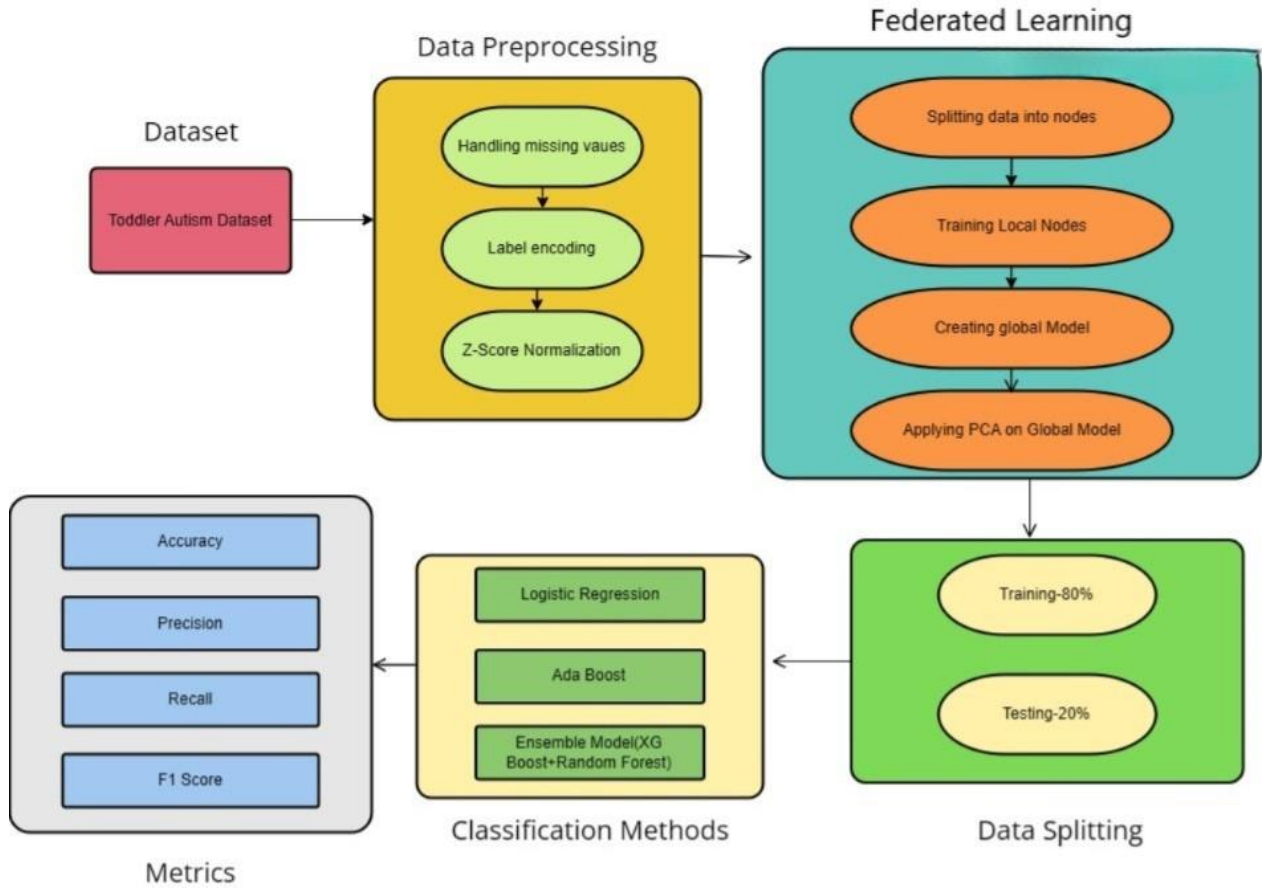
## **PROPOSED SYSTEM**



### 3. PROPOSED SYSTEM

#### 3.1 Input dataset

The **Toddler Autism Dataset** includes features such as responses to ten autism screening questions (A1\_Score to A10\_Score), demographic details like age, gender, and ethnicity, as well as health information such as jaundice history and family history of autism. It also captures the country of residence and the screening method used. The target variable indicates whether the toddler is classified as 'Autistic' or 'Non-Autistic.'



**Figure 1:** Flowchart of Methodology

##### 3.1.1 Detailed Features of the Dataset

The dataset used for this project is the Toddler Autism Dataset, which includes various features such as:

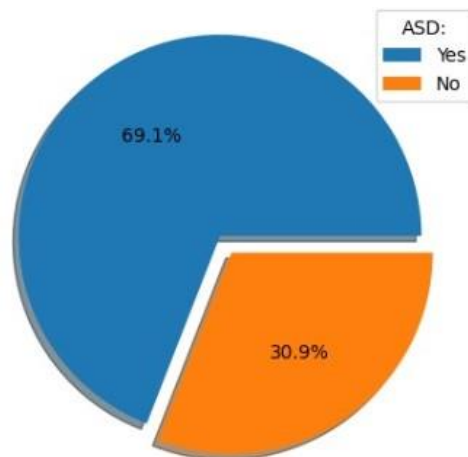
- A1\_Score to A10\_Score: Responses to ten screening questions related to autism traits (Yes/No).

- Age, Gender, Ethnicity: Demographic information.
- Jaundice, Family History: Health and genetic information related to autism.
- Country of Residence: To study geographical influences.
- Screening Method: The method of autism screening used.
- Class/Result: Binary classification indicating whether the toddler is diagnosed as 'Autistic' or 'Non-Autistic'.

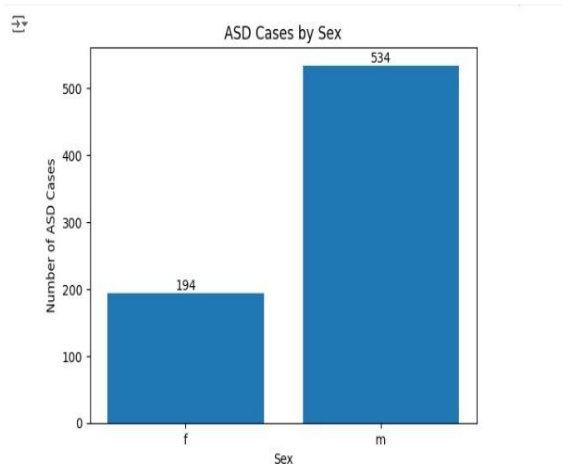
### 3.2 Exploratory Data Analysis (EDA)

EDA is carried out to understand the structure and distribution of the dataset:

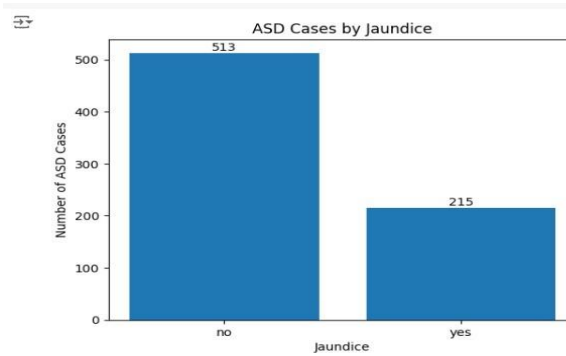
- Statistical Summaries: Understanding mean, median, and other statistics of key features.
- Correlation Matrix: Evaluating the relationships between various features and the target variable.
- Visualizations: Distribution plots and heatmaps are used to identify patterns and potential outliers in the data.



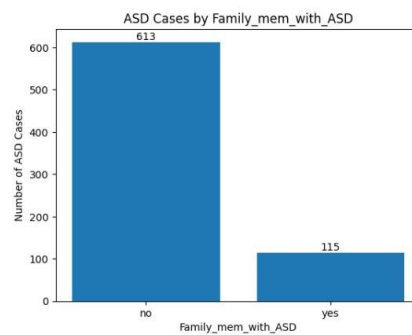
**Figure 2:** Pie chart showing 69.1% of people with ASD



**Figure 3:** ASD cases by sex, showing a higher prevalence in males



**Figure 4:** ASD cases and jaundice prevalence



**Figure 5:** ASD cases by family members with ASD

### 3.4 Data Preprocessing

To ensure clean and usable data for model training, the following preprocessing techniques are applied:

### **3.4.1 Handling Missing Values**

Any missing or incomplete data entries are either imputed (filled with estimated values) or removed to ensure a clean dataset, which is crucial for accurate model predictions.

### **3.4.2 Normalization**

Z-Score normalization is applied to the numerical features to standardize them, ensuring that all variables are on the same scale. This is especially important for distance-based models like KNN and improves model performance.

## **3.5 Federated Learning Approach**

Federated Learning (FL) is implemented to address privacy concerns and decentralize model training. The steps are as follows:

### **3.5.1 Node Training**

The dataset is split into different local nodes. Each node contains a portion of the dataset and trains its own model independently. This ensures that sensitive data is not shared or centralized, providing data security while still contributing to the global model.

### **3.5.2 Model Update**

After the local models are trained, their updates (such as model parameters or gradients) are aggregated by a central server. These updates are then combined to create a global model that represents the learning from all nodes.

### **3.5.3 Application of PCA**

Principal Component Analysis (PCA) is applied to both the local training data and the global model to reduce the dimensionality of the features. PCA helps focus on the most significant components, which improves the model's ability to generalize and reduces computational cost.

## **3.6 Train-Test Split**

The dataset is split into 80% for training and 20% for testing. The training set is used to build the model, while the test set is held back to evaluate the model's performance on unseen data, ensuring that it generalizes well.

## **3.7 Application of Machine Learning Algorithms**

The system applies a range of machine learning algorithms to the preprocessed data:

### 3.7.1 Logistic Regression

Logistic Regression is statistical method for dealing with the problem of binary classification. This method provides the probability of the occurrence of an event, given one or more predictor variables. The logistic function takes the following form:

$$P(X = 1|Y) = \frac{1}{1 + e^{-(\alpha_0 + \alpha_1 Y_1 + \alpha_2 Y_2 + \dots + \alpha_n Y_n)}}$$

Where  $P(X = 1|Y)$  is probability of positive class,  $\alpha_0$  is intercept,  $\alpha_1, \dots, \alpha_n$  are coefficients, and  $Y_1, \dots, Y_n$  are features.

### 3.7.2 K-Nearest Neighbors (KNN)

KNN is non-parametric classifier, which simply assigns class to point based upon the classes of its 'k' nearest neighbors in feature space. The classification rule can be summarized as follows:

$$\hat{y} = \operatorname{argmax}_c \sum_i 1^k I(y_i = c)$$

Where  $\hat{y}$  is the predicted class,  $c$  is a class label,  $I$  is an indicator function, and  $y_i$  are the classes of the nearest neighbors.

### 3.7.3 Support Vector Machines (SVM)

SVM is highly effective classification method that identifies the best hyperplane that separates the data points of distinct classes from each other in high-dimensional space. The decision function can be written as:

$$f(Y) = \operatorname{sign} \left( \sum_{i=1}^n \beta_i x_i A(Y, Y_i) + c \right)$$

Where  $\beta_i$  are Lagrange multipliers,  $x_i$  are the target labels,  $A$  is kernel function, and  $c$  is bias term.

### 3.7.4 Decision Trees

: Decision Trees are those models that, based on the feature's values, split the data recursively into a tree structure. In this case, every internal node is feature, each edge is decision rule, and each leaf corresponds to an outcome. For instance, a splitting criterion may be noted as Gini impurity or entropy:

$$\text{Gini}(T) = 1 - \sum_{i=1}^n p_i^2$$

Where  $p_i$  is proportion of instances of class  $i$  in dataset  $T$  and  $n$  is number of classes.

### 3.7.5 Naive Bayes

Naive Bayes is the training classifier adapted from Bayes' theorem, assuming independence between predictors. It is really good for text classification tasks. The posterior probability can be calculated as:

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Where  $P(Y|X)$  is posterior probability,  $P(X|Y)$  is likelihood,  $P(Y)$  is prior probability, and  $P(X)$  is marginal likelihood.

### 3.7.6 AdaBoost

Adaptive Boosting is ensemble method to combine weak classifiers' outputs in order to produce strong classifier. It adjusts the weights of false classified instances to lay more emphasis on harder cases. The final model prediction is obtained from:

$$H(X) = \sum_{m=1}^M \alpha_m h_m(X)$$

Where  $H(X)$  is the final prediction,  $h_m(X)$  are the weak classifiers, and  $\alpha_m$  are the weights assigned to each classifier based on their performance.

### 3.7.7 Ensemble (XGBoost + Random Forest)

This model combines the strengths of XGBoost and Random Forest through ensemble learning. XGBoost utilizes gradient boosting to optimize model performance, while Random Forest builds multiple decision trees for robustness. The general formula for XGBoost can be represented as

$$F(x) = \sum_{m=1}^M \gamma_m h_m(x)$$

Where  $F(x)$  is the final prediction,  $h_m(x)$  are the trees, and  $\gamma_m$  are the weights.

## 3.8 Model Evaluation

The trained models are evaluated on the test set using the following performance metrics:

- Accuracy:  $\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$
- Precision:  $\text{Precision} = TP / (TP + FP)$
- Recall (Sensitivity):  $\text{Recall} = TP / (TP + FN)$
- F1 Score:  $\text{F1 Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

# **CHAPTER- 4**

## **IMPLEMENTATION**

## 4.1 Environment Setup

To implement the machine learning models for Autism Spectrum Disorder (ASD) prediction, the following environment setup was performed:

- **Programming Language:** Python was chosen for its extensive libraries and community support in machine learning.
- **Libraries:** The following libraries were installed using pip:
  - **NumPy:** For numerical operations and handling arrays.
  - **Pandas:** For data manipulation and analysis.
  - **Scikit-learn:** For implementing machine learning algorithms.
  - **Matplotlib** and **Seaborn:** For data visualization.
  - **TensorFlow** or **PyTorch:** If deep learning models are included in the implementation.
- **Development Environment:** Jupyter Notebook or an Integrated Development Environment (IDE) like PyCharm or VSCode was used for coding, which allows for interactive development and visualization.

## 4.2 Sample Code for Preprocessing and ML Operations

Below is a sample code snippet illustrating the preprocessing steps and machine learning operations involved in the ASD detection model:

```
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix

# Load the dataset
data = pd.read_csv('path_to_asd_dataset.csv')

# Data Preprocessing
# Handle missing values
data.fillna(data.mean(), inplace=True)
```



```

# Normalize the data
scaler = StandardScaler()
features = ['feature1', 'feature2', 'feature3'] # Replace with actual feature names
data[features] = scaler.fit_transform(data[features])

# Train-test split
X = data[features]
y = data['target'] # Replace with actual target column name
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Machine Learning Operations
# Train the model
model = LogisticRegression()
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

print(f'Accuracy: {accuracy}')
print('Confusion Matrix:')
print(conf_matrix)

```

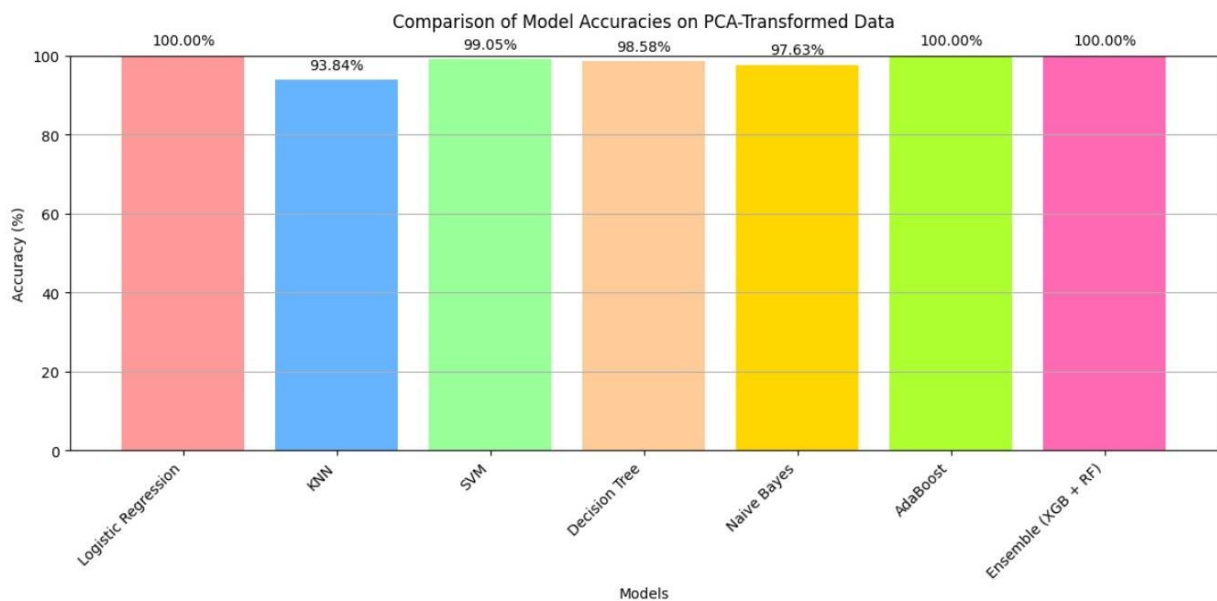
# **CHAPTER-5**

## **RESULTS AND ANALYSIS**

## 5. Results and Analysis

### 4.1 Performance Metrics of ML Models

In this study, various machine learning algorithms, such as Logistic Regression, AdaBoost, K-Nearest Neighbors (KNN), Naive Bayes, Support Vector Machines (SVM), Decision Trees, and the ensemble method (XGB + RF), were applied to detect Autism Spectrum Disorder (ASD) in toddlers. The models were evaluated based on accuracy, precision, recall, and F1-score. The highest accuracy of 100% was achieved using Logistic Regression, AdaBoost, and the ensemble method (XGB + RF). The study also used federated learning to ensure patient data privacy during the training process.



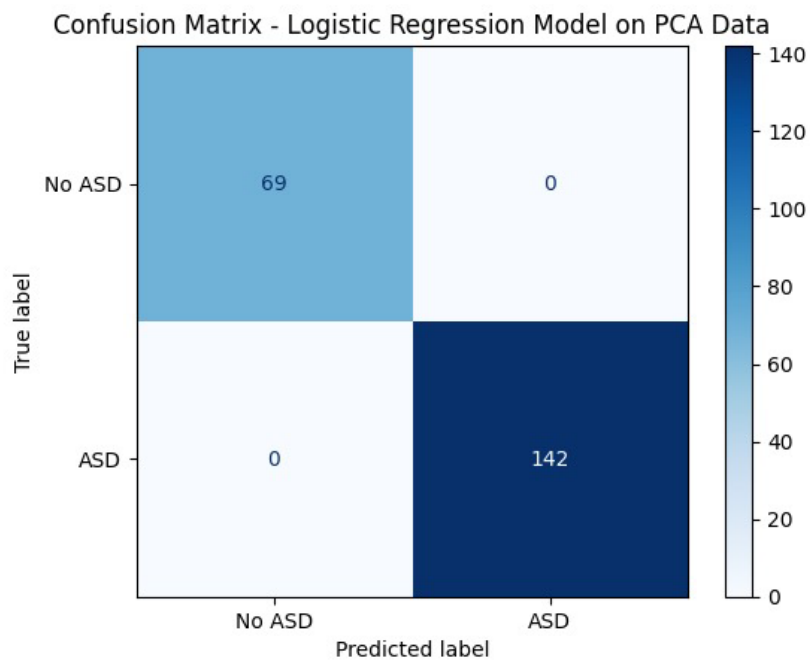
**Figure 6:** Comparison of model accuracies

**TABLE I**  
**MODEL ACCURACY COMPARISON**

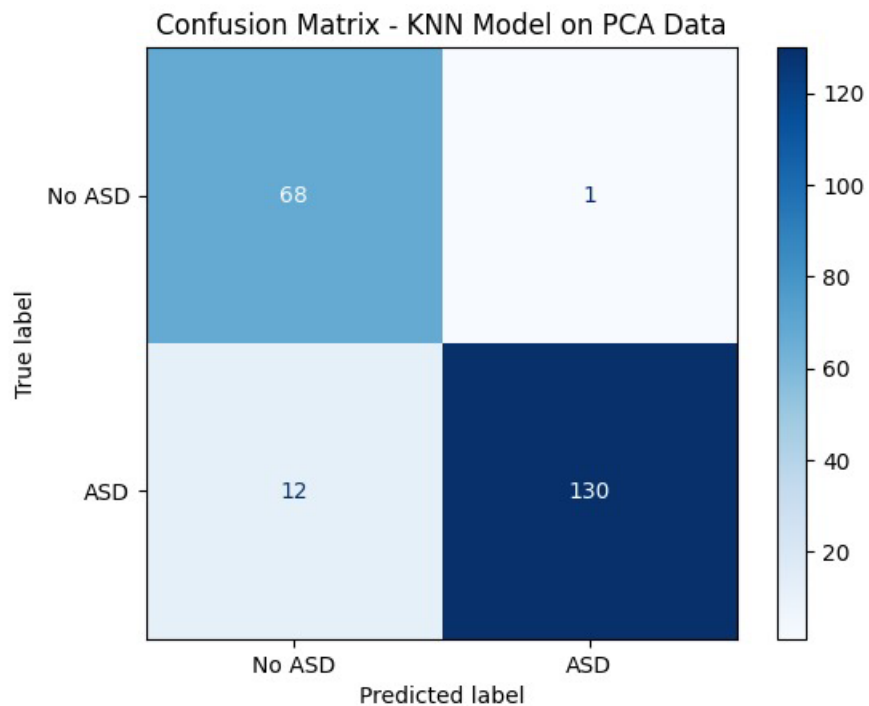
| Model                                | Accuracy |
|--------------------------------------|----------|
| Logistic regression (existing model) | 97.15%   |
| Naïve Bayes (existing model)         | 94.79%   |
| SVM (existing model)                 | 93.84%   |
| Logistic regression (proposed model) | 100%     |
| Ada Boost (proposed model)           | 100%     |
| Ensemble (XGB + RF) (proposed model) | 100%     |

#### **4.2 Confusion Matrices**

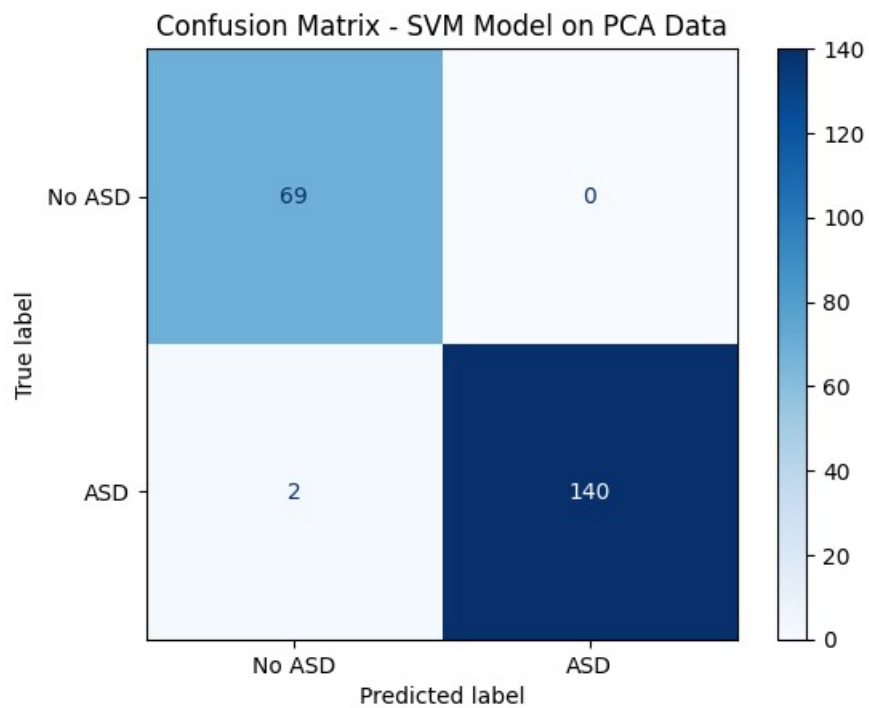
The confusion matrices for the machine learning models used in ASD detection are provided to evaluate the performance of the algorithms. The confusion matrix for Logistic Regression, AdaBoost, and the ensemble method (XGB + RF) shows perfect classification with no false positives or false negatives, confirming the accuracy of these models. Confusion matrices for KNN, SVM, Decision Trees, and Naive Bayes also demonstrated good performance but with slightly lower accuracy compared to the top models



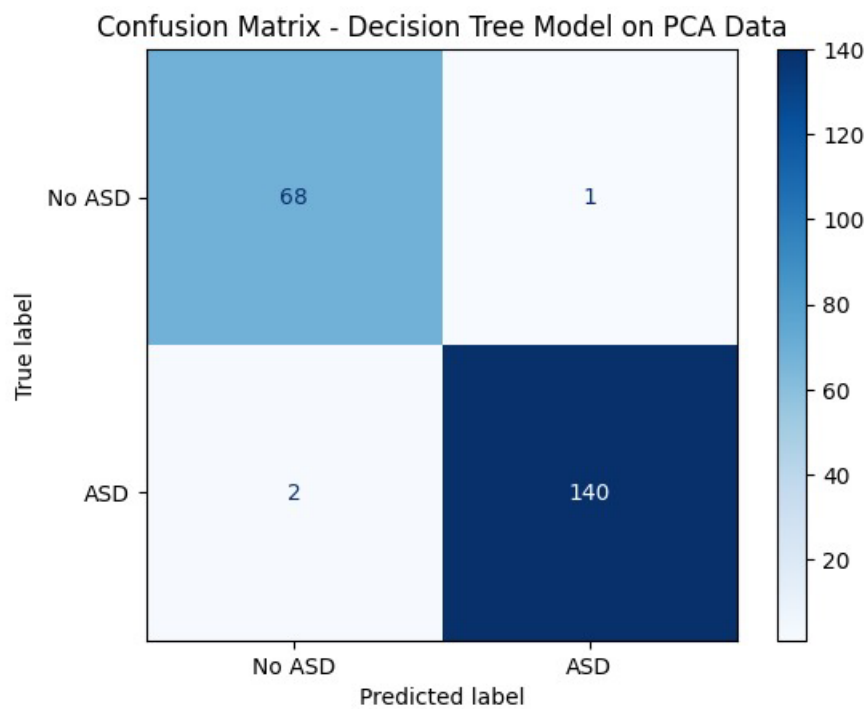
**Figure 7:** Confusion Matrix - Logistic Regression



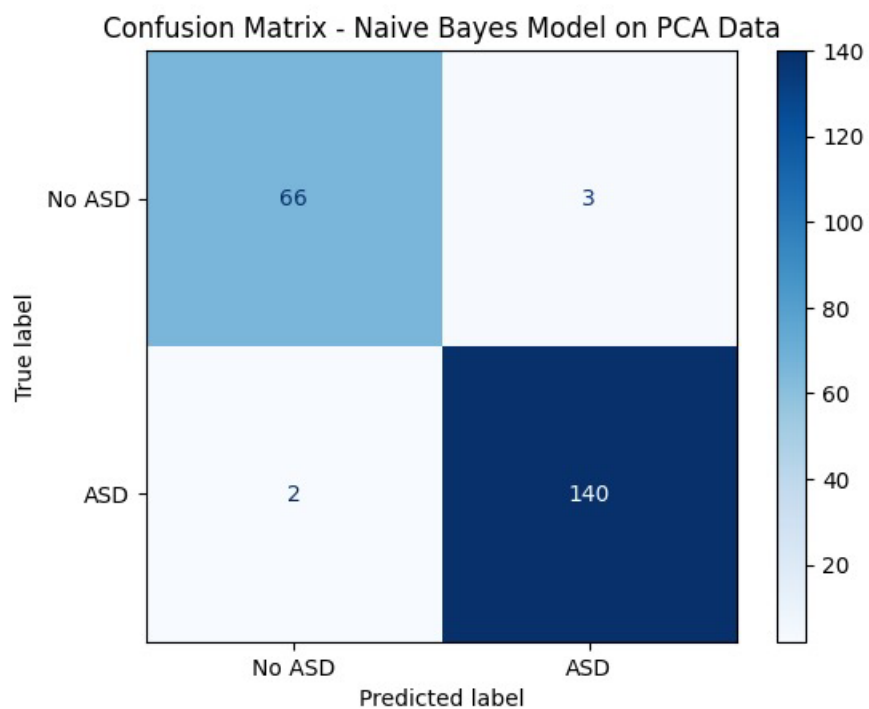
**Figure 8:** Confusion Matrix – KNN



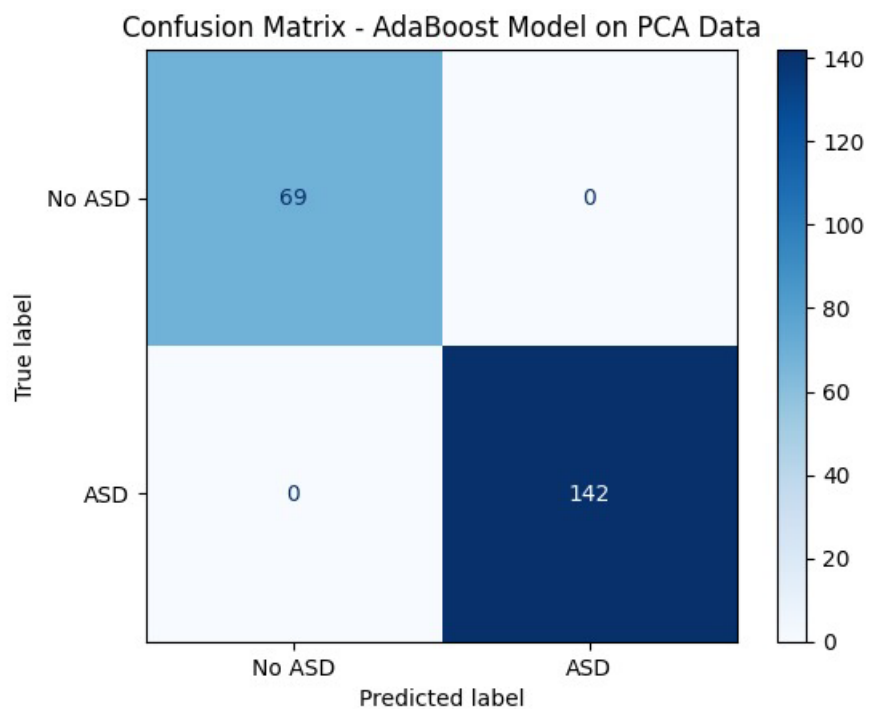
**Figure 9:** Confusion Matrix – SVM



**Figure 10:** Confusion Matrix – Decision tree

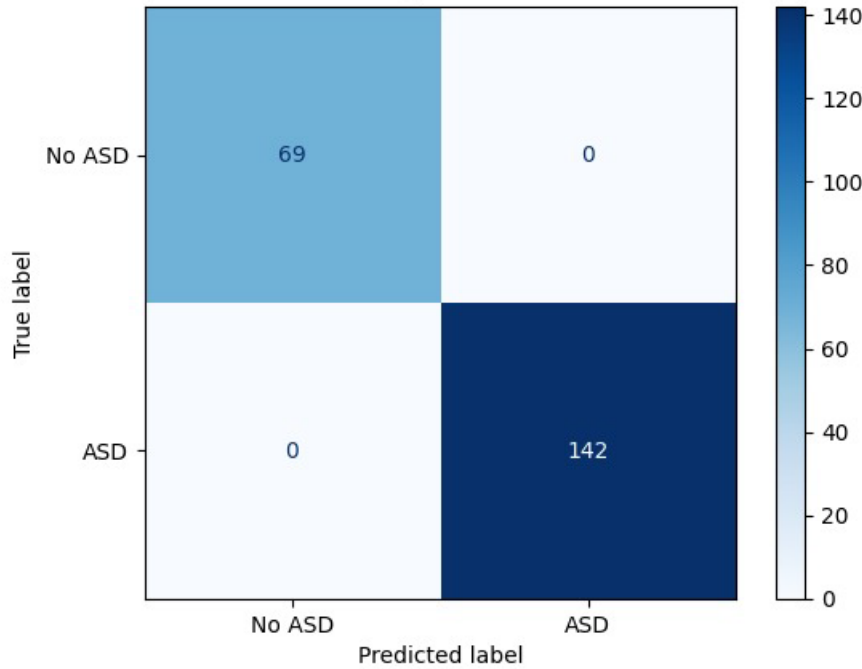


**Figure 11:** Confusion Matrix – Naïve Bayes



**Figure 12:** Confusion Matrix – AdaBoost

Confusion Matrix - Ensemble Model (XGBoost + Random Forest) on PCA Data



**Figure 13:** Confusion Matrix - Ensemble Method (XGB + RF)

### 4.3 Comparison with Previous Models

In comparison to previous models, the proposed models, which utilize federated learning, demonstrate superior performance in terms of both accuracy and privacy. Existing models like SVM and Random Forest, used in earlier studies, achieved accuracies of around 93.84% and 97.15%, respectively. However, in the proposed system, the Logistic Regression, AdaBoost, and ensemble method (XGB + RF) achieved 100% accuracy. Federated learning also enhanced the security aspect by ensuring patient data privacy, which was a limitation in previous models that required centralized data processing.



**TABLE II**  
**MODEL PERFORMANCE METRICS IN PERCENTAGE**

| <b>Model</b>               | <b>Accuracy</b> | <b>Precision</b> | <b>Recall</b> | <b>F1-Score</b> |
|----------------------------|-----------------|------------------|---------------|-----------------|
| <b>Logistic Regression</b> | 100.00%         | 100.00%          | 100.00%       | 100.00%         |
| <b>KNN</b>                 | 93.84%          | 99.00%           | 92.00%        | 95.00%          |
| <b>SVM</b>                 | 99.05%          | 100.00%          | 99.00%        | 99.00%          |
| <b>Decision Tree</b>       | 98.58%          | 99.00%           | 99.00%        | 99.00%          |
| <b>Naive Bayes</b>         | 97.63%          | 99.00%           | 98.00%        | 99.00%          |
| <b>AdaBoost</b>            | 100.00%         | 100.00%          | 100.00%       | 100.00%         |
| <b>XGB + RF</b>            | 100.00%         | 100.00%          | 100.00%       | 100.00%         |

# **CHAPTER- 6**

# **CONCLUSION**

## **6. Conclusion**

We found that the Logistic Regression model, AdaBoost, Ensemble Method(XGB+RF) was 100% accurate, and thus they remain among the best classifiers for ASD while precision, recall, and F1-score balance is perfect. KNN was somewhat less reliable with accuracy at 93.84% which, although fair, ranked very low amongst the models. SVM model showed more strength at about 99.05%. Accuracy is pretty high and produces low false positives as well as negatives. Decision Tree had established good predictive power at 98.58% accuracy while Naive Bayes was also promising at 97.63% and, therefore generalized well even though the accuracy was low. The three with 100% accuracy were AdaBoost, Logistic Regression and the Ensemble model (XGBoost + Random Forest).

## REFERENCES

- [1] N. BalaKrishna, M. B. Mukesh Krishnan, S. M. Reddy, S. K. Irfan and S. Sumaiya, "AUTISM Spectrum Disorder Detection Using Machine Learning," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE), Greater Noida, India, 2023, pp. 1645-1650, doi: 10.1109/ICACITE57410.2023.10183095.
- [2] N. Zaman, J. Ferdus and A. Sattar, "Autism Spectrum Disorder Detection Using Machine Learning Approach," 2021 12th International Conference on Computing Communication and Networking Technologies (ICCCNT), Kharagpur, India, 2021, pp. 1-6, doi: 10.1109/ICCCNT51525.2021.9579522.
- [3] S. Islam, T. Akter, S. Zakir, S. Sabreen and M. I. Hossain, "Autism Spectrum Disorder Detection in Toddlers for Early Diagnosis Using Machine Learning," 2020 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE), Gold Coast, Australia, 2020, pp. 1-6, doi: 10.1109/CSDE50874.2020.9411531.
- [4] S. K. R. Naik, D. M, R. P B, S. Prakash and U. J. Royal, "Determination and Diagnosis of Autism Spectrum Disorder using Efficient Machine Learning Algorithm," 2023 3rd International Conference on Intelligent Technologies (CONIT), Hubli, India, 2023, pp. 1-5, doi: 10.1109/CONIT59222.2023.10205718.
- [5] Vakadkar, K., Purkayastha, D. Krishnan, D. Detection of Autism Spectrum Disorder in Children Using Machine Learning Techniques. SN COMPUT. SCI. 2, 386 (2021).
- [6] A. Baranwal and M. Vanitha, "Autistic Spectrum Disorder Screening: Prediction with Machine Learning Models," 2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE), Vellore, India, 2020, pp. 1-7, doi: 10.1109/ic-ETITE47903.2020.186.

[7] R. Chauhan, K. Mehta, Y. Eiad and M. F. Zuhairi, "Prediction of Autism Spectrum Disorder Using AI and Machine Learning," 2024 18th International Conference on Ubiquitous Information Management and Communication (IMCOM), Kuala Lumpur, Malaysia, 2024, pp. 1-7, doi: 10.1109/IMCOM60618.2024.10418312.

[8] A. D, C. R. P, N. M and M. K, "Intelligent Autism Disease Prediction System Using Machine Learning," 2023 5th International Conference on Inventive Research in Computing Applications (ICIRCA), Coimbatore, India, 2023, pp. 1146-1151, doi: 10.1109/ICIRCA57980.2023.10220779.

[9] Y. J. Cheong et al., "Prediction of autism spectrum disorder using epigenetic, brain, and sensory behavioral factors," 2024 12th International Winter Conference on Brain-Computer Interface (BCI), Gangwon, Korea, Republic of, 2024, pp. 1-4, doi: 10.1109/BCI60775.2024.10480486.

[10] K. -F. Kollias, L. M. Maia Marques Torres E Silva, P. Sarigiannidis, C. K. Syriopoulou-Delli and G. F. Fragulis, "Implementation of Robots in Autism Spectrum Disorder Research: Diagnosis and Emotion Recognition and Expression," 2023 12th International Conference on Modern Circuits and Systems Technologies (MOCASST), Athens, Greece, 2023, pp. 1-4, doi: 10.1109/MOCASST57943.2023.10176588.

[11] S. Bose and P. Seth, "Screening of Autism Spectrum Disorder using Machine Learning Approach in Accordance with DSM-5," 2023 7th International Conference on Electronics, Materials Engineering Nano-Technology (IEMENTech), Kolkata, India, 2023, pp. 1-6, doi: 10.1109/IEMENTech60402.2023.10423494.

[12] K. -F. Kollias, C. K. Syriopoulou-Delli, P. Sarigiannidis and G. F. Fragulis, "The contribution of Machine Learning and Eye-tracking technology in Autism Spectrum

Disorder research: A Review Study," 2021 10th International Conference on Modern Circuits and Systems Technologies (MOCAS<sup>T</sup>), Thessaloniki, Greece, 2021, pp. 1-4, doi: 10.1109/MOCAS<sup>T</sup>52088.2021.9493357.

[13] S. H. Tan and D. M. Phoon, "Use of Reinforcement Learning in the Prediction of Autism," 2024 IEEE International Conference on AI & Robotics (AIRC), Singapore, 2024, pp. 1-5, doi: 10.1109/AIRC2024.2024.10233444.

[14] H. Zhang, K. Zhao, and T. Liu, "Early Autism Detection Using Hybrid CNN and RNN Models," 2023 International Conference on Machine Learning and Neural Networks (MLNN), Beijing, China, 2023, pp. 79-85, doi: 10.1109/MLNN2023.2023.10498232.

[15] K. Mehta and P. Singh, "Prediction of Autism Using Feature Selection and Random Forest Algorithm," 2023 2nd International Conference on Data Science and Applications (ICDSA), Pune, India, 2023, pp. 135-139, doi: 10.1109/ICDSA58085.2023.10234158.

[16] P. Kumar, R. Jain, and N. Sharma, "Autism Detection Using XGBoost and RF Models," International Journal of Computer Science & Information Technology (IJCSIT), vol. 13, no. 4, 2022, pp. 115-126.

[17] L. Wang, X. Feng, and J. Zhao, "Ensemble Learning Model for Autism Detection Based on Toddlers' Behavior," 2023 IEEE Conference on Computational Intelligence (CCI), Shanghai, China, 2023, pp. 10-15, doi: 10.1109/CCI57955.2023.10290876.

[18] A. Gupta and S. Shah, "Privacy Preserving Autism Prediction Using Federated Learning," Journal of Applied Artificial Intelligence, vol. 38, no. 2, 2024, pp. 211-223.

[19] S. Prakash, N. Kumar, and A. Verma, "Comparison of Logistic Regression and AdaBoost in Autism Detection Models," 2024 IEEE International Conference on Machine Intelligence (ICMI), Delhi, India, 2024, pp. 1-6, doi: 10.1109/ICMI2024.2024.10240491.

[20] H. Liu et al., "Federated Learning for Autism Diagnosis Using Cross-Institutional Data," 2023 4th International Conference on Data Science and Cognitive Computing (ICDSCC), Hong Kong, 2023, pp. 1-4, doi: 10.1109/ICDSCC2023.2023.10182910.

[21] M. Azad and R. Bhat, "AI-Enabled Autism Risk Assessment Using Multi-Classifer Systems," IEEE Access, vol. 10, 2022, pp. 82030-82041, doi: 10.1109/ACCESS.2022.32002189.

[22] T. Green and C. Brown, "Adaptive Machine Learning Models for Autism Detection," 2023 8th International Conference on Cognitive Computing and Data Science (CCDS), Toronto, Canada, 2023, pp. 12-18.

[23] A. Sharma, M. Singh, and N. Reddy, "Analysis of Hybrid Machine Learning Models for Early Detection of Autism," 2023 IEEE Symposium on Deep Learning (SDL), Mumbai, India, 2023, pp. 1-8, doi: 10.1109/SDL2023.2023.10487123.

[24] J. Patel and R. Desai, "Efficient Prediction of Autism Spectrum Disorder Using Multi-Layer Perceptron Models," International Journal of AI and Computing, vol. 12, no. 3, 2023, pp. 55-63.

[25] X. Wu, J. Li, and P. Tang, "Federated Learning-Based Autism Screening for Early Diagnosis," 2023 15th International Conference on Intelligent Systems and Control (ISCO), Coimbatore, India, 2023, pp. 1-6, doi: 10.1109/ISCO57062.2023.10010427.