

Early stage detection of Autism Spectrum Disorder using AI and Machine Learning

Name of Authors: - Sujal Ghodke, Vedant Joshi, Suyash Patil, Dhananjay Gawali

ABSTRACT

It is always a complex procedure to diagnosis autism spectrum disorder (ASD) because there is no particular medical test for autism, like a blood test, to make diagnosis for the disorder. Autism spectrum disorder is define by the disability and impairments of participating in social communication or the presence of restricted or repetitive behaviors, or both. It is a behaviorally diagnosed condition. To make a diagnosis, doctor look at the child's developmental history and behavior. Apparently, most children do not attain a proper diagnosis for autism until it is too late. Sometimes parents find it difficult to accept that their child's cerebral development is not keeping pace with their physical development. This lateness in diagnosis hinders a child's ability to get the help they need to keep developing. It is important to diagnose ASD as early as possible through monitoring, screening, and reach their full potential. So, we are developing a system that will have the ability to diagnosis Autism and come up with reliable and effective conclusion even without the help of professional. We hope this system will be very helpful for those concerned parents who are worried about their child's growth and activities, at the time same time it will be very useful for the professionals.

INTRODUCTION

Autism spectrum disorder (ASD) is a complicated neurological condition that impairs communication, behaviour, and social interaction. The prevalence of ASD has been steadily increasing worldwide, making it a significant public health concern. Early diagnosis and intervention are crucial for optimizing the developmental trajectory of individuals with ASD.

While there have been advancements in ASD diagnosis, current methods can be time-consuming and subjective. Machine learning, a subset of artificial intelligence, offers a promising approach to enhance the accuracy and efficiency of early ASD detection. By analyzing large amounts of data, machine learning models can identify subtle patterns and correlations that may not be apparent to human experts.

This study aims to develop a machine learning-based system for the early detection of ASD. By utilizing a combination of clinical and behavioral data, the proposed system aims to provide a more accurate and efficient diagnostic tool. Through the application of advanced machine learning techniques, we seek to improve the early detection of ASD, enabling timely intervention and support for individuals with this condition.

NEED OF THE STUDY.

The study of Autism Spectrum Disorder (ASD) when using machine learning is necessary due to the ability to bring revolution into preliminary identity, diagnosis and treatment. Machine learning algorithm can analyze behavior patterns, developmental piles and various biological markers, enabling time intervention, to detect ASD at an early stage, which improves the results of the development. Traditional clinical methods are often subjective, doctors chair too much on observations and parents' reports. However, machine learning provides a more purpose and data -driven approach by checking the pattern in speech, eye movement and social interactions. These models can facilitate mass screening, especially in areas with limited access to health services through mobile applications and online tools. In addition, machine learning ASD helps identify sub -factories, provides insight into the diverse nature of the disorder and supports individual treatment plans. By analyzing the results of the treatment, these models can guess which intervention can do the best work for a person and improve medical efficiency. In addition, machine learning technology helps researchers highlight genetic, neurological and environmental factors that contribute to ASD, which provide deep insight into its biological basis. Integration of machine learning into ASD research thus strongly promises to continue our understanding, diagnosis and treatment of disorders, eventually improves the quality of life for individuals on spectrum and their families.

3.1 Population and Sample

In the context of Autism Spectrum Disorder (ASD), a population refers to the complete organization of folks that may be studied to apprehend, diagnose, or deal with the disorder. For instance, the population may want to encompass all youngsters international or within a selected location who're liable to ASD. Studying this complete population is frequently impractical because of the large number of people and the assets required. Therefore, researchers use a sample, that is a smaller, consultant group of individuals decided on from the populace. For example, a examine would possibly examine behavioral patterns in a pattern of one,000 children from numerous demographics to identify early indicators of ASD. This sample, if selected successfully, can provide valuable insights that can be generalized to the larger population. Machine mastering models can then use this sample information to discover patterns, expect hazard elements, and guide early diagnosis, in the end improving the accuracy and performance of ASD research and interventions.

3.2 Data and Sources of Data

The take a look at of Autism Spectrum Disorder (ASD) using device learning has won huge attention, with researchers exploring techniques like guide vector machines, selection trees, and logistic regression to enhance prognosis accuracy (Thabtah, 2017). Thabtah, Kamalov, and Rajab (2018) proposed a computational intelligence framework to enhance screening, at the same time as Vaishali and Sasikala (2018) used binary firefly optimization to achieve 92.12% accuracy. Early prognosis is essential, as Zwaigenbaum et al. (2015) highlighted the advantages of early intervention in improving cognitive and social capabilities. Genetic and environmental elements additionally contribute to ASD, with Lyall et al. (2017) noting that eighty% of cases are inherited, alongside risks from pollutants and infections. Challenges persist because of ASD's heterogeneity, moral concerns in genetic research, and diagnostic complexities. Future studies targets to integrate AI and ML for actual-time monitoring and personalized interventions, at the same time as addressing social and conversation problems, inclusive of struggles with understanding others' mental states (Baron-Cohen et al., 1985).

3.3 Theoretical framework

The theoretical framework of Autism Spectrum Disorder (ASD) encompasses various views to apprehend its reasons, characteristics, and developmental styles. Neurobiological theories propose that ASD is linked to structural and practical abnormalities in brain areas like the amygdala, hippocampus, and prefrontal cortex, with bizarre neural connectivity and neurotransmitter imbalances, specifically in serotonin and dopamine. Genetic and epigenetic theories highlight the hereditary nature of ASD, with about eighty% of cases being inherited, while environmental elements for the duration of prenatal improvement can also have an impact on gene expression. Cognitive theories, along with the Theory of Mind (Baron-Cohen et al., 1985), endorse that individuals with ASD struggle to interpret others' mind and feelings, whilst govt disorder and weak primary coherence theories provide an explanation for demanding situations with planning, flexibility, and contextual understanding. Behavioral theories, particularly those primarily based on Applied Behavior Analysis (ABA), emphasize how reinforcement and getting to know can form behavior, contributing to effective interventions. Environmental elements, such as prenatal publicity to pollution and superior parental age, are also identified as potential ASD danger factors. Additionally, system mastering fashions are now being applied to ASD studies, assisting to analyze behavioral and genetic information for early diagnosis and intervention. This multidisciplinary framework offers a complete understanding of ASD, helping in the improvement of extra powerful diagnostic equipment and customized remedy techniques.

RESEARCH METHODOLOGY:-

This look at appears at younger children, between 18 and 24 months vintage, who would possibly display early signs of autism. Finding out if a baby has autism early is surely important because it allows them get the right aid. The researchers used a collection of information from the UCI Machine gaining knowledge of Repository, which has facts from 1,050 youngsters-some of whom are traditional (702 children) and some who have autism (348 children). They amassed 21 exclusive pieces of records using unique exams, like the changed tick list for autism in infants (M-CHAT). To make certain the observe consists of an amazing mic of children, they cautiously selected a number of ages, genders, and ranges of signs. Since the statistics became already accumulated an doesn't have names on it, they didn't want unique permission to apply it.

3.1 Population and Sample

Population: The target population includes children aged 2–12 years at risk of Autism Spectrum Disorder (ASD), encompassing diverse demographics (gender, ethnicity, socioeconomic status) across urban and rural settings. Sample:

- Sample Size: A representative sample of 1,500 children, split into 70% training and 30% testing datasets.
- Inclusion Criteria: Children exhibiting early ASD markers (e.g., delayed speech, limited eye contact) or parental/clinical concerns.
- Exclusion Criteria: Children with other neurodevelopmental disorders (e.g., Down syndrome) or severe sensory impairments.
- Sampling Technique: Stratified random sampling to ensure diversity in age, gender, and geographic location.
- Ethical Considerations: Informed consent from guardians, anonymization of data, and approval from institutional review boards

3.2 Data and Sources of Data

Data Types:

- 1. Clinical Data: Diagnostic evaluations (ADOS-2, ADI-R), medical history, and genetic markers.
- 2. **Behavioural Data:** Eye-tracking metrics, social interaction patterns, and repetitive behaviour logs.
- 3. **Demographic Data:** Age, gender, parental age, and socioeconomic status.

Sources:

- **Primary Data:** Collected via partnerships with paediatric clinics and ASD screening centres.
- Secondary Data: Public datasets (e.g., UCI Autism Screening Dataset, NDAR, ABIDE) and peer-reviewed studies.
- Tools: Structured questionnaires, mobile apps for parental reporting, and wearable sensors for real-time monitoring.

Preprocessing:

- Missing data handled via k-nearest neighbours (KNN) imputation.
- Feature scaling (normalization) and one-hot encoding for categorical variables.
- Class imbalance addressed using Synthetic Minority Oversampling Technique (SMOTE)

In this look at, researchers used 3 main forms of statistics to construct and check device learning fashions: scientific facts, behavioral facts, and demographic records. Clinical data included diagnostic exams, scientific histories, and genetic facts to better recognize autism. Behavioral records centered on interest styles, social interactions, and repetitive behaviors, while demographic information protected age, gender, and own family heritage to pick out danger factors. The information was carefully cleaned, missing values had been expected the use of comparable information, and express variables had been transformed into numbers for version compatibility. To stability the dataset, techniques were applied to create extra autism case examples. The statistics came from primary sources like kids's clinics and autism screening facilities, and secondary resources like public databases and beyond research research, ensuring a diverse and dependable dataset for model training.

3.3 Theoretical framework

The theoretical foundation of this study is built upon a multidisciplinary integration of computational intelligence, neurobiological insights, cognitive-behavioral theories, and device mastering principles, together designed to cope with the complexities of Autism Spectrum Disorder (ASD) prognosis. At its center, the framework leverages neurobiological theory, which posits that unusual mind connectivity and neurotransmitter imbalances underpin ASD's behavioral manifestations. This concept informs the incorporation of physiological facts, consisting of eye-monitoring metrics and electroencephalogram (EEG) recordings, to capture deviations in interest styles and neural synchrony. For example, eye-tracking records quantify gaze fixation during social interactions, reflecting deficits in joint interest—an indicator of ASD—at the same time as EEG indicators screen aberrant oscillatory interest in mind areas related to social cognition, together with the prefrontal cortex and amygdala. Complementing this, cognitive concept, specially the Theory of Mind (ToM), emphasizes the demanding situations people with ASD face in attributing intellectual states to others. This theoretical lens guides the analysis of social interaction metrics, consisting of responsiveness to emotional cues or reciprocity in dialogue, that are quantified through dependent observational coding and wearable sensor information. These metrics operationalize To M deficits, translating abstract cognitive impairments into measurable variables for gadget gaining knowledge of models.

Further enriching the framework is behavioral idea, rooted in Applied Behavior Analysis (ABA), which examines how reinforcement and environmental interactions form repetitive or restrictive behaviors. This concept directs the selection of features which include frequency of stereotyped circulate

3.4 Statistical tools and econometric models

This section elaborates the proper statistical/econometric/financial models which are being used to forward the study from data towards inferences. The detail of methodology is given as follows.

3.4.1 Data Preprocessing and Analysis Tools

The study employed Python's scikit-learn library for data manipulation, model implementation, and evaluation. Pandas was utilized for dataset cleaning, including handling missing values through median imputation and encoding categorical variables (e.g., gender, symptom severity) into numerical formats. To address the inherent class imbalance (ASD:non-ASD \approx 1:2), the Synthetic Minority Oversampling Technique (SMOTE) was applied, synthetically generating ASD cases to balance the dataset. Feature normalization was performed using Z-score standardization to ensure uniformity in scale across variables such as social interaction scores and repetitive behavior metrics.

3.4.2 Feature Selection and Dimensionality Reduction

The Binary Firefly Algorithm (BFA), a swarm intelligence technique, was implemented to identify the most predictive features from the UCI dataset. This metaheuristic approach simulates firefly behavior, where "brighter" fireflies (features with higher diagnostic relevance) attract others, iteratively optimizing the feature subset. BFA reduced the original 21 features to a critical subset of 10, including response to name, eye contact duration, and sensory sensitivities, which collectively explained 85% of variance in ASD prediction. This aligns with Vaishali et al.'s (2018) methodology, ensuring computational efficiency and minimizing overfitting.

3.4.3 Machine Learning Classifiers

- Support Vector Machines (SVM) with a radial basis function (RBF) kernel were used to handle non-linear decision boundaries. The regularization parameter C and kernel coefficient gamma were tuned via grid search to maximize separability between ASD and non-ASD classes.
- Decision Trees were constructed using the CART algorithm, with pruning (max depth = 5) to prevent overfitting. Feature importance was calculated using Gini impurity reduction.

Logistic Regression with L2 regularization (ridge regression) served as the baseline model, assessing linear relationships between features and ASD likelihood

3.4.4 Evaluation Metrics and Validation

Model performance was quantified using accuracy, precision, recall, F1-score, and ROC-AUC. To ensure robustness, 10-fold cross-validation was applied, partitioning the dataset into training and validation subsets iteratively. Stratified sampling preserved the class distribution across folds, mitigating bias. Additionally, confusion matrices were generated to visualize false positives/negatives, while SHAP (Shapley Additive explanations) values provided interpretability by quantifying feature contributions to predictions.

3.4.5 Econometric Adjustments for Bias Mitigation

To address biases from imbalanced data and parent-reported behavioral assessments, stratified sampling weights were incorporated during model training, ensuring proportional representation of minority (ASD) cases. Probabilities from SMOTE-augmented datasets were calibrated using Platt scaling to align predicted likelihoods with true ASD prevalence rates. Furthermore, bootstrap resampling (1,000 iterations) was employed to estimate confidence intervals for accuracy metrics, enhancing the reliability of results in real-world applications.

3.4.6 Software and Visualization Tools

Matplotlib and Seaborn libraries were used to visualize feature distributions, ROC curves, and decision boundaries. TensorFlow integration enabled preliminary exploration of neural networks, though traditional classifiers outperformed deep learning models due to limited dataset size. Jupyter Notebooks facilitated reproducible workflows, documenting every stage from preprocessing to final model deployment.

Decision Trees were constructed using the CART algorithm, with pruning (max depth = 5) to prevent overfitting. Feature importance was calculated using Gini impurity reduction.

Logistic Regression with L2 regularization (ridge regression) served as the baseline model, assessing linear relationships between features and ASD likelihood.

IV. RESULTS AND DISCUSSION

4.1 Results of Descriptive Statics of Study Variables

Model	Accuracy	Error rate	Precision	Recall	F1-score	AUC-ROC
SVM	93.4	6.6	0.91	0.94	0.92	0.96
Decision Tree	89.1	10.9	0.87	0.82	0.84	0.88
Logistic Regression	87.5	12.5	0.85	0.83	0.84	0.89

Key Findings:

The proposed machine studying gadget for early ASD detection demonstrated robust performance, making use of the Binary Firefly Algorithm (BFA) for function optimization and SMOTE for sophistication imbalance mitigation. BFA decided on 10 key behavioral functions from the original 21 in the UCI dataset, with top predictors like reaction to call, eye touch length, and repetitive motor actions, together explaining eighty five% of the variance in ASD prediction. The hybrid SVM-BFA version outperformed different classifiers, attaining 93.Four% accuracy and a 0.96 AUC-ROC, surpassing previous research consisting of Vaishali et al. (2018). Decision Trees confirmed overfitting notwithstanding pruning, while Logistic Regression exhibited better specificity on the price of sensitivity. SMOTE drastically stepped forward version sensitivity, reducing false negatives from 22% to 9%, making sure better early diagnosis—a vital issue in ASD screening. Clinically, the version aligns with DSM-5 standards, with functions like eye touch period and repetitive behaviors validating its relevance. However, challenges persist, which include do not forget bias in figure-suggested statistics and restrained dataset size. Future improvements should contain integrating real-time information from wearables, incorporating multimodal biomarkers like neuroimaging, and developing cell applications to increase accessibility, specially in low-aid settings.

I. ACKNOWLEDGMENT

We would love to increase our sincere gratitude to anybody who supported us during this venture. First and primary, we're immensely thankful to our assignment guide, Dr. Rupesh Mahajan, for their valuable steerage, encouragement, and knowledge, which had been instrumental within the a success of entirety of this paintings. We also increase our heartfelt way to the Head of the Department, Dr. Vinod Kimbahune, for presenting an inspiring and supportive mastering surroundings at Dr. D.Y. Patil Institute of Technology. Their management and willpower to academic excellence had been a steady source of motivation. Finally, we desire to express our private appreciation to our own family, friends, and everybody who contributed without delay or circuitously to the a hit crowning glory of this project.

REFERENCES

- 1. K. M. Dalton, B. M. Nacewicz, T. Johnstone et al., "Gaze fixation and the neural circuitry of face processing in autism;" *Nature Neuroscience*, vol. 8, no. 4, pp. 519-526, 2005.
- 2. Kim M. Dalton, B. M. Nacewicz, A. L. Alexander, and R. J. Davidson, "Gaze-fixation, brain activation, and amygdala volume in unaffected siblings of individuals with autism;' *Biological Psychiatry*, vol. 61, no. 4, pp. 512-520, 2007.
- 3. C. Ecker, S. Y. Bookheimer, and D. G. Murphy, "Neuroimaging in autism spectrum disorder: brain structure and function across the lifespan;' *The Lancet Neurology*, vol. 14, no. 11,
- 4. pp. 1121-1134,2015.
- 5. M. Plitt, K. A. Barnes, and A. Martin, "Functional connectivity classification of autism identifies highly predictive brain features but falls short of biomarker standards;' *Neuroimage: Clinical*, vol. 7, pp. 359-366, 2015.
- 6. B. Mwangi, K. P. Ebmeier, K. Matthews, and J. Douglas Steele, "Multi-centre diagnostic classification of individual structural neuroimaging scans from patients with major depressive disorder; '*Brain*, vol. 135, no. 5, pp. 1508-1521, 2012.
- 7. 0. Tadevosyan-Leyfer, M. Dowd, R. Mankoski et al., 'J\. principal components analysis of the autism diagnostic interview revised;' *Journal of the American Academy of Child & Adolescent Psychiatry*, vol. 42, no. 7, pp. 864-872, 2003.
- 8. K. ermand, W A. Ghani, and A. I. Shihab, "Classification and monitoring of autism using SVM and VMCM;" *Journal of Theoretical and Applied Information Technology*, vol. 96, no. 14, 2018.
- 9. L. Cattaneo, M. Fabbri-Destro, S. Boria et al., "Impairment of actions chains in autism and its possible role in intention understanding;' *Proceedings of the National Academy of Sciences*, vol. 104, no. 45, pp. 17825-17830, 2007.
- 10. A. I. Shihab, "Data classification and applied bioinformatics for monitoring of autism using neural network;' *Journal of Engineering and Applied Sciences*, vol. 13, no. 5 SI, pp. 4778-4785, 2018.

