# Analysis Report on Machine Learning for Game-Based ASD Detection

## 1. Abstract:

Autism Spectrum Disorder (ASD) is a neurodevelopmental disorder characterized by difficulties in social interaction, communication, and repetitive behaviors. Early detection and diagnosis are essential for timely intervention, which can significantly improve outcomes. Traditional diagnostic approaches often involve structured assessments and interviews, which can be time-consuming and stressful for children.

This report explores the use of machine learning algorithms to analyze behavioral data collected through game-based environments, aiming to identify early signs of ASD. By applying algorithms such as Logistic Regression, k-Nearest Neighbors (KNN), Decision Trees, and Neural Networks, we can detect patterns in user behavior that may indicate ASD. The focus is on using simple, interpretable models and neural networks to classify behaviors based on response times, social interactions, and decision-making patterns without relying on computer vision techniques.

## 2. Problem Statement:

The key challenge addressed in this project is the early detection of Autism Spectrum Disorder (ASD) in children through game-based behavioral assessments. We aim to create a machine learning-based model that can:
1. Accurately classify children as ASD-positive or ASD-negative based on their in-game behaviors.
2. Identify patterns such as delayed response to social cues, repetitive actions, and atypical decision-making strategies that are commonly associated with ASD.
3. Provide an interpretable, accessible solution that clinicians can use as a supplementary diagnostic tool alongside traditional methods.

The primary problem lies in extracting meaningful features from gameplay data and applying machine learning models that can generalize well to different children, balancing accuracy with interpretability and ease of deployment.

## 3. Dataset Consideration:

The dataset will consist of behavioral data collected during game interactions. Some typical features extracted from the gameplay might include:
- Response Time: Time taken to respond to social or task-related cues.
- Interaction Frequency: Frequency of interactions with in-game characters or tasks.
- Repetitive Actions: Count of repeated behaviors (e.g., performing the same task multiple times).
- Decision-Making: Pattern of choices made in social scenarios (e.g., helping others, making socially appropriate responses).
- Social Engagement: How often the child engages in social behaviors like communication or cooperation.

The dataset will have both ASD-positive and ASD-negative labels, with features representing behavioral patterns and target labels indicating the diagnosis.

## 4. Possible Machine Learning Models:

In this project, we focus on traditional machine learning algorithms and neural networks. Below is a summary of models we will consider, along with their advantages and challenges.

### A. Logistic Regression (LR):

**Description**:
Logistic Regression is a simple, interpretable model that estimates the probability of a binary outcome. For ASD detection, it can be used to predict whether the child falls on the autism spectrum based on the behavioral features extracted from gameplay.

- **Advantages**:
  - Simple, interpretable, and easy to implement.
  - Provides probabilities, making it useful for threshold-based decisions.
  - Performs well on linearly separable data.

- **Challenges**:
  - Limited to linear decision boundaries, which may not capture complex behaviors associated with ASD.
  - Sensitive to outliers.

- **When to Use**:
  - Useful as a baseline model for binary classification.
  - Appropriate when feature relationships with ASD labels are approximately linear.

### B. k-Nearest Neighbors (KNN):

**Description**:
KNN is a non-parametric, instance-based learning algorithm that classifies a data point based on the majority label of its k nearest neighbors. It can be applied to classify behaviors into ASD-positive or ASD-negative groups based on proximity to similar behavior patterns.

- **Advantages**:
  - Simple to understand and implement.
  - Non-linear: Can model more complex relationships between features.
  - Works well with smaller datasets.

- **Challenges**:
  - Performance degrades with high-dimensional data (many features).
  - Sensitive to noisy data and irrelevant features.
  - Computationally expensive as the dataset grows.

- **When to Use**:
  - Suitable when the data has clear local patterns or clusters.
  - Effective for small datasets with meaningful proximity relations.

## C. Decision Trees:

**Description**:
Decision Trees use a tree-like model of decisions and their possible consequences. For ASD detection, they can be used to model decision-making pathways based on game behavior, identifying key features (e.g., repetitive actions, response delays) that lead to classification.

- **Advantages**:
  - Highly interpretable; decisions are easy to understand.
  - Can handle both numerical and categorical data.
  - No need for feature scaling.
  - Handles non-linear relationships well.

- **Challenges**:
  - Prone to overfitting, especially with complex trees.
  - Can be unstable (small changes in data may result in a different tree).

- **When to Use**:
  - When interpretability is a priority (clinicians may want to understand which behaviors lead to ASD classification).
  - When non-linear relationships between features are expected.

## D. Random Forest:

**Description**:
Random Forest is an ensemble learning method that builds multiple decision trees and merges them to improve accuracy and avoid overfitting. It can provide insights into which game behaviors are most indicative of ASD.

- **Advantages**:
  - Reduces overfitting compared to a single decision tree.
  - Handles high-dimensional data well.
  - Provides feature importance scores, helping to identify the most significant behaviors.

- **Challenges**:
  - Less interpretable than a single decision tree.
  - Can be computationally intensive with large datasets.

- **When to Use:**
  - When the dataset is large and overfitting is a concern.
  - When we want to understand which features (behaviors) are most important in determining ASD.

## E. Support Vector Machines (SVM):

**Description**:
SVM constructs a hyperplane or set of hyperplanes in a high-dimensional space to separate the data into different classes. For ASD detection, it can help classify children based on non-linear relationships between behavioral features.

- **Advantages**:
  - Effective in high-dimensional spaces.

- Can model non-linear boundaries using kernels.
- Robust to outliers.

- **Challenges**:
  - Computationally expensive with large datasets.
  - Harder to interpret compared to simpler models.

- **When to Use**:
  - When the dataset is relatively small, and complex, non-linear decision boundaries are needed.

### F. Neural Networks (NN):

**Description**:
Neural networks are powerful models capable of learning complex, non-linear relationships between features and labels. A simple feedforward neural network can be used to classify ASD-positive and ASD-negative children based on their game behavior.

- **Advantages**:
  - Can model very complex patterns in data.
  - Scalable to larger datasets and high-dimensional features.

- **Challenges**:
  - Requires large datasets for good performance.
  - Harder to interpret compared to traditional models like Decision Trees or Logistic Regression.
  - Prone to overfitting if not properly regularized.

- When to Use:
  - When the dataset is large, and simpler models fail to capture the complexity of behavioral patterns.
  - When the goal is high accuracy, even at the cost of interpretability.

## 5. Model Selection and Evaluation:

We will begin with simpler models (e.g., Logistic Regression, KNN, Decision Trees) to establish baseline performance. These models are interpretable and easy to deploy, which is valuable when the goal is clinical use.

Once baseline models are implemented, more complex models such as Random Forest and Neural Networks will be considered. The following metrics will be used to evaluate model performance:
- **Accuracy**: Overall correctness of the model.
- **Precision and Recall**: Important for understanding how well the model identifies true ASD-positive cases (recall) and how often predicted positives are correct (precision).
- **F1 Score**: A balanced measure of precision and recall.
- **AUC-ROC Curve**: To evaluate the model's ability to distinguish between ASD-positive and ASD-negative classes.

Cross-validation techniques (e.g., k-fold cross-validation) will be employed to ensure that the models generalize well to unseen data.