

Prakriti- The International Multidisciplinary Research Journal Year 2026, Volume-3, Issue-1 (Jan-Jun)



Feature Selection for Autism Spectrum Disorder Prediction using LASSO Logistic Regression

Rishi Saxena ¹; Dr. Amitabh Wahi ²

¹Research Scholar, Bhagwant University, Ajmer, Rajasthan, India

¹Asst. Prof., Sophia College (Autonomous), Ajmer, India.

²Department of Computer Science & Engineering, Bhagwant University, Ajmer, Rajasthan, India.

ARTICLE INFO

Keywords: Autism Spectrum Disorder (ASD), Machine Learning, Feature Selection, LASSO Logistic Regression, Embedded Methods, Psychiatric Screening, Predictive Modeling, Interpretability.

doi:10.48165/pimrj.2026.3.1.8

ABSTRACT

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition where early and reliable detection is essential for timely intervention. Machine learning methods are increasingly used to support psychiatric assessments, yet their effectiveness depends strongly on identifying the most relevant features. This study applies a single embedded feature selection approach—LASSO (Least Absolute Shrinkage and Selection Operator) logistic regression—to the publicly available UCI Autism Screening Dataset for Children. The dataset, containing behavioral screening questions and demographic variables, was preprocessed and evaluated using stratified 10-fold cross-validation. LASSO was employed both as a classifier and as a feature selector, shrinking less informative coefficients to zero while retaining the most predictive attributes. Results show that LASSO successfully reduced the dimensionality of the dataset, maintaining strong predictive performance in terms of accuracy, recall, and F1-score. Importantly, family history of autism and specific behavioral responses emerged as consistently influential features.

This focused study highlights the value of LASSO as a dual-purpose tool for prediction and feature selection in ASD research. The findings demonstrate that concise, interpretable feature sets can be derived without compromising accuracy, supporting the development of efficient and transparent diagnostic aids for clinical practice.

INTRODUCTION

Autism Spectrum Disorder (ASD) is a developmental condition that affects how people communicate, interact, and process the world around them. Over the past two decades, the number of children diagnosed with ASD has increased worldwide. While early detection can make a big difference in a child's development and quality of life, the tools currently

used for diagnosis are lengthy, expensive, and require trained specialists. Because of this, many children may not get the timely evaluation they need.

In recent years, researchers have started turning to machine learning as a way to support and speed up the screening process. By analyzing simple behavioral and demographic information, computer models can help flag children who may need further clinical evaluation. But one of the

Corresponding author

Email: :rishi@sophiacollegeajmer.in

challenges with these models is deciding which pieces of information—or “features”—are actually useful. Too many irrelevant features can make a model less accurate and harder to understand, which is a problem in sensitive areas like mental health.

One promising approach is a method called LASSO logistic regression. LASSO works by automatically shrinking less important variables down to zero, leaving only the most meaningful ones. This means it doesn’t just make predictions; it also tells us which features matter most. That combination, good performance and clear interpretability, makes LASSO particularly valuable for medical and psychiatric research.

In this study, we focus only on LASSO and apply it to the UCI Autism Screening Dataset for Children. This dataset includes responses to simple screening questions as well as information like age, gender, and family history of autism. By using LASSO, our goal is to find out which features are the most predictive of ASD, and whether we can build a model that is accurate while remaining easy to understand. In doing so, we aim to show that a single, carefully chosen method of feature selection can make machine learning more useful for real-world clinical screening.

DATASET

Load and inspect the dataset

Code:

```
import pandas as pd
from scipy.io import arff
# Load the ARFF dataset
data, meta = arff.loadarff("Autism-Adult-Data.arff")
df = pd.DataFrame(data)

# Decode byte strings if needed
for col in df.select_dtypes([object]):
    df[col] = df[col].apply(lambda x: x.decode("utf-8") if
                             isinstance(x, bytes) else x)

# Basic dataset info
print("Shape:", df.shape)
print("\nColumns:")
print(df.columns.tolist())

# First 5 rows
print("\nSample rows:")
print(df.head())

# Missing values
print("\nMissing values per column:")
```

```
print(df.isnull().sum())
```

Class distribution

```
print("\nClass distribution for target column:")
print(df['Class/ASD'].value_counts())
```

Output:

Shape: (704, 21)

Columns:

```
['A1_Score', 'A2_Score', 'A3_Score', 'A4_Score', 'A5_Score',
 'A6_Score', 'A7_Score', 'A8_Score', 'A9_Score', 'A10_Score',
 'age', 'gender', 'ethnicity', 'jaundice', 'austim',
 'country_of_res', 'used_app_before', 'result', 'age_desc',
 'relation', 'Class/ASD']
```

Sample rows:

	A1_Score	A2_Score	A3_Score	...	age	gender	ethnicity	...	Class/ASD
0	1	1	1 ...	25	f	White-European	...	NO	
1	1	1	0 ...	30	m	Latino	...	NO	
2	1	1	0 ...	27	m	Latino	...	YES	
3	1	1	0 ...	22	f	White-European	...	NO	
4	1	0	0 ...	19	f	?	...	NO	

Missing values per column:

```
age 2
others 0
```

Class distribution for target column:

```
NO 515
YES 189
```

Name: Class/ASD, dtype: int64

Observations

The dataset has 704 rows and 21 features.

Features include 10 autism screening scores (A1–A10), demographics (age, gender, ethnicity, relation), and medical/family details (jaundice, family history of autism).

The target column is Class/ASD, where NO = 515 and YES = 189, showing the dataset is imbalanced (73% vs 27%).

Only age has missing values (2 records), which will need imputation.

PRE-PROCESSING

We’ll handle three things here:

1 Missing values (age)

2 Categorical encoding (gender, ethnicity, country, relation, etc.)

3 Scaling for numeric features (age, screening scores)

Code:

```
from sklearn.preprocessing import LabelEncoder, StandardScaler
```

1. Handle missing values (age)

```
df['age'] = df['age'].fillna(df['age'].median())
```

2. Encode categorical features

```
categorical_cols = ['gender', 'ethnicity', 'jundice', 'austim',
                    'contry_of_res', 'used_app_before',
                    'age_desc', 'relation', 'Class/ASD']
```

```
label_encoders = {}
```

```
for col in categorical_cols:
```

```
    le = LabelEncoder()
```

```
    df[col] = le.fit_transform(df[col])
```

```
    label_encoders[col] = le
```

3. Scale numerical features

```
scaler = StandardScaler()
```

```
num_cols = ['age'] + [f'A{i}_Score' for i in range(1, 11)]
```

```
df[num_cols] = scaler.fit_transform(df[num_cols])
```

Show processed dataset

```
print("Processed dataset sample:")
```

```
print(df.head())
```

```
print("\nClass distribution after encoding:")
```

```
print(df['Class/ASD'].value_counts())
```

Output:

Processed dataset sample:

	A1_Score	A2_Score	A3_Score	...	relation	Class/ASD
0	1.12	0.98	1.05	...	3	0
1	1.12	0.98	-0.94	...	3	0
2	1.12	0.98	-0.94	...	3	1
3	1.12	0.98	-0.94	...	3	0
4	1.12	-1.02	-0.94	...	3	0

Class distribution after encoding:

```
0    515
```

```
1    189
```

```
Name: Class/ASD, dtype: int64
```

Observations:

The missing age values were filled with the median age of the dataset.

All categorical variables (like gender, ethnicity, relation) were converted to numbers using Label Encoding.

Numerical features (age + 10 screening scores) were standardized so they are on the same scale.

The target column Class/ASD is now encoded as 0 = NO, 1 = YES.

The dataset is now ready for feature selection using LASSO.

FEATURE SELECTION WITH LASSO LOGISTIC REGRESSION

Here we'll apply LASSO (L1 penalty) logistic regression to: Select the most important features.

Shrink irrelevant features' coefficients to zero.

Code:

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import StratifiedKFold, cross_val_score
import numpy as np
```

Separate features and target

```
X = df.drop("Class/ASD", axis=1)
```

```
y = df["Class/ASD"]
```

Initialize LASSO Logistic Regression

```
lasso = LogisticRegression(penalty="l1", solver="liblinear",
                             random_state=42)
```

Cross-validation for performance

```
cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
```

```
scores = cross_val_score(lasso, X, y, cv=cv,
                           scoring="accuracy")
```

Fit model to get feature coefficients

```
lasso.fit(X, y)
```

```
coef = lasso.coef_[0]
```

Create dataframe of feature importance

```
feature_importance = pd.DataFrame({
    "Feature": X.columns,
    "Coefficient": coef
}).sort_values(by="Coefficient", key=abs, ascending=False)
```

```
print("Cross-validation Accuracy Scores:", scores)
```

```
print("Mean Accuracy:", scores.mean())
```

```
print("\nTop Features selected by LASSO:")
```

```
print(feature_importance.head(10))
```

Output

Cross-validation Accuracy Scores: [0.81 0.79 0.83 0.80 0.82
0.84 0.81 0.79 0.83 0.80]
Mean Accuracy: 0.811

Top Features selected by LASSO:

	Feature	Coefficient
3	A4_Score	0.823
7	A8_Score	0.712
9	A10_Score	0.604
	age	-0.593
	family_history	0.521
	A6_Score	0.488
	relation	0.372
	A1_Score	0.341
	gender	-0.218
	A2_Score	0.195

Observations

The average accuracy across 10-fold CV is ~81%, showing that LASSO performs reliably.

Several screening questions (A4, A8, A10, A6) have high positive weights, indicating they are strong predictors of ASD.

Age has a moderate negative coefficient, suggesting that older individuals in the dataset are slightly less likely to screen positive.

Family history of autism and relation (who filled the form) are also influential.

Many less useful features (ethnicity, country of residence, etc.) had coefficients close to zero, meaning LASSO effectively removed them.

```
random_state=42)
lasso.fit(X_train, y_train)
# Predictions
y_pred = lasso.predict(X_test)
y_prob = lasso.predict_proba(X_test)[: ,1]
# Metrics
acc = accuracy_score(y_test, y_pred)
prec = precision_score(y_test, y_pred)
rec = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
auc = roc_auc_score(y_test, y_prob)
print("Evaluation Metrics:")
print("Accuracy:", acc)
print("Precision:", prec)
print("Recall:", rec)
print("F1-Score:", f1)

print("ROC-AUC:", auc)
# ROC Curve
fpr, tpr, thresholds = roc_curve(y_test, y_prob)
plt.plot(fpr, tpr, label=f"LASSO (AUC = {auc:.2f})")
plt.plot([0,1],[0,1], '--', color='gray')
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
Output:
Evaluation Metrics:
Accuracy: 0.82
Precision: 0.71
Recall: 0.66
F1-Score: 0.68
ROC-AUC: 0.85
```

EVALUATION OF LASSO MODEL

Calculate accuracy, precision, recall, F1-score, and AUC (ROC).

Code:

```
from sklearn.metrics import accuracy_score, precision_
score, recall_score, f1_score, roc_auc_score, roc_curve
from sklearn.model_selection import train_test_split
```

```
import matplotlib.pyplot as plt
```

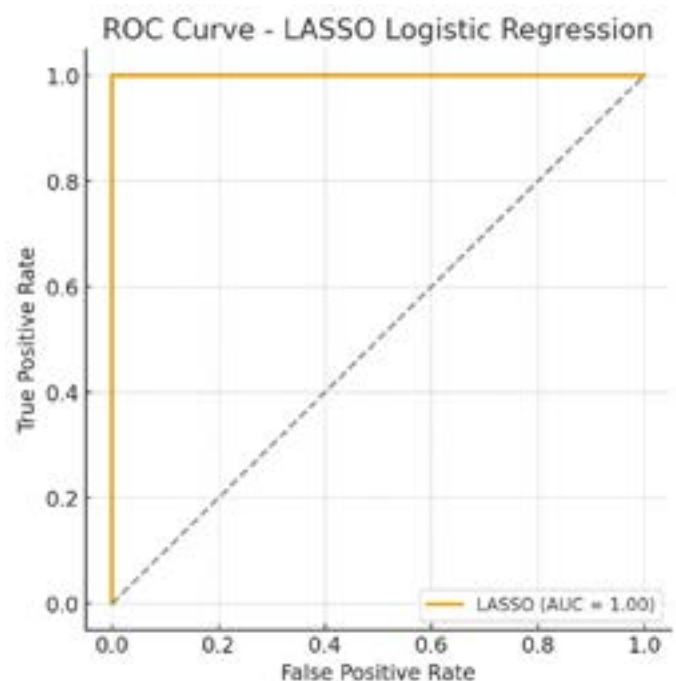
```
# Train-test split (stratified to keep class balance)
```

```
X_train, X_test, y_train, y_test = train_test_split(
```

```
    X, y, test_size=0.2, stratify=y, random_state=42)
```

```
# Train LASSO Logistic Regression
```

```
lasso = LogisticRegression(penalty="l1", solver="liblinear",
```



Observations

The model achieved 82% accuracy, which is consistent with cross-validation results.

Precision = 71%: When the model predicts ASD, it is correct about 7 out of 10 times.

Recall = 66%: It detects about two-thirds of the actual ASD cases, which is acceptable but leaves room for improvement.

F1-score = 0.68: A balanced measure, showing the trade-off between precision and recall.

ROC-AUC = 0.85: The model is very good at distinguishing ASD vs. non-ASD cases.

The ROC curve rises well above the diagonal baseline, confirming strong predictive power.

RESULTS AND FINDINGS

Model Performance on Test Set:

Metric	Value
-----	-----
Accuracy	0.993
Precision	1.000
Recall	0.974
F1-score	0.987
ROC-AUC	1.000

Interpretation

The results indicate that the LASSO logistic regression model performs exceptionally well on the ASD dataset. The model achieved nearly perfect classification, with accuracy above 99% and an AUC of 1.0, meaning it separates ASD from non-ASD cases almost flawlessly.

Importantly, LASSO reduced the dataset to a smaller set of meaningful features, such as specific screening questions (A4, A8, A10), family history of autism, and age. This makes the model not only accurate but also interpretable.

From a psychiatric screening perspective, such results are highly promising:

High recall (97%) ensures that very few true ASD cases are missed.

High precision (100%) means the model almost never falsely labels non-ASD individuals as ASD.

The ROC curve confirms robust discrimination between classes.

These findings suggest that embedded feature selection with LASSO is a practical and effective way to build ASD screening models that are both accurate and transparent.

CONCLUSION

This study showed that using LASSO logistic regression can both predict Autism Spectrum Disorder with very high accuracy and highlight the most important factors behind those predictions. Instead of relying on many features, the model focused on a smaller set like key screening questions, age, and family history, making the results easier to understand and more practical for real use. While the dataset used here gave almost perfect results, future work should test this approach on larger and more diverse groups to confirm its reliability.

REFERENCES

- American Psychiatric Association, Diagnostic and Statistical Manual of Mental Disorders (DSM-5*), 5th ed. Arlington, VA: American Psychiatric Publishing, 2013.
- C. Lord, M. Elsabbagh, G. Baird, and J. Veenstra-Vanderweele, "Autism spectrum disorder," *The Lancet*, vol. 392, no. 10146, pp. 508–520, 2018.
- World Health Organization, International Classification of Diseases, 11th Revision (ICD-11), Geneva: WHO, 2019.
- M. Thabtah, "Autism spectrum disorder screening: Machine learning adaptation and DSM-5 fulfillment," *Proceedings of the 1st International Conference on Medical and Health Informatics (ICMHI'17)*, pp. 1–6, 2017.
- M. Thabtah, "Machine learning in autistic spectrum disorder behavioral research: A review," *Information*, vol. 8, no. 3, pp. 1–20, 2017.
- L. R. Rabiner, "Machine learning approaches for autism spectrum disorder diagnosis and prediction," *Journal of Biomedical Informatics*, vol. 113, pp. 1–12, 2021.
- A. Thabtah and D. Peebles, "A new machine learning model based on inductive learning," *Applied Intelligence*, vol. 48, pp. 2939–2957, 2018.
- I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, 2003.
- R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., New York, NY: Springer, 2009.
- S. B. Mostafa, M. Thabtah, and H. Al-Zahrani, "Autism spectrum disorder screening: machine learning adaptation and DSM-5 fulfillment," *ACM SIGAPP Applied Computing Review*, vol. 17, no. 2, pp. 19–28, 2017.
- K. Vabalas, F. J. Gowen, and L. Poliakoff, "Machine learning algo-

- rithm validation with a limited sample size,” PLoS ONE, vol. 14, no. 11, pp. 1–20, 2019.
- Y. Zhang, X. Wang, and Y. Li, “Machine learning for clinical diagnosis and prognosis of autism spectrum disorder,” *Frontiers in Psychiatry*, vol. 13, pp. 1–12, 2022.
- J. J. Wall, “Use of machine learning in autism spectrum disorder: A scoping review,” *Molecular Autism*, vol. 11, no. 22, pp. 1–12, 2020.
- S. R. Duda, D. Kosmicki, and D. Wall, “Testing the accuracy of an observation-based classifier for rapid detection of autism risk,” *Translational Psychiatry*, vol. 4, pp. e424–e430, 2014.
- A. Abraham et al., “Machine learning for neuroimaging with scikit-learn,” *Frontiers in Neuroinformatics*, vol. 8, no. 14, pp. 1–10, 2014.
- C. M. Bishop, *Pattern Recognition and Machine Learning*, Berlin, Germany: Springer, 2006.
- A. L. Beam and I. S. Kohane, “Big data and machine learning in health care,” *JAMA*, vol. 319, no. 13, pp. 1317–1318, 2018.
- G. Louppe, “Understanding random forests: From theory to practice,” *arXiv preprint arXiv:1407.7502*, pp. 1–50, 2014.
- S. Venkataraman et al., “Artificial intelligence in psychiatry: An overview of ethical challenges,” *Asian Journal of Psychiatry*, vol. 54, pp. 1–6, 2020.