

# Research Log - ML SEM-4 PROJECT

---

Avantika Agarwal

TOPIC: Multimodal ML models integrating Genomic and medical information for ASD detection

# EXPERIMENT-1

**DATE: 23/01/24**

## Model Description and Input

The model used was Logistic Regression and the input was a Breast cancer dataset present as a default dataset in scikit learn ( for practice purposes). In this model, the probabilities describing the possible outcomes of a single trial are modeled using a logistic function. In this code, we are carrying out Binary Logistic Regression, meaning, there would be two possible outcomes of the prediction ( that is, it predicts the presence or absence of breast cancer - 0/1 ).

## Results

PARAMETER	EXPLANATION	VALUE
Accuracy score (normalized)	The accuracy_score function computes the accuracy→ the fraction (default) in this case	0.951048951048951
Accuracy score (no normalization)	The accuracy_score function computes the accuracy→ the count (normalize=False) of correct predictions.	136
Area under the curve (AUC)	AUC - ROC curve is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve and AUC represents the degree or measure of separability. It tells how much the model is capable of distinguishing between classes.  Compute Area Under the Curve (AUC) using the trapezoidal rule.	0.9937106918238994

	This is a general function, given points on a curve	
Sensitivity	Sensitivity = $TP / (TP + FN)$	0.9444444444444444 4
Specificity	Specificity = $TN / (TN + FP)$	0.962264150943396 2
Precision	The precision is the ratio $tp / (tp + fp)$ where $tp$ is the number of true positives and $fp$ the number of false positives. The precision is intuitively the ability of the classifier not to label as positive a sample that is negative.  The best value is 1 and the worst value is 0.	0.977011494252873 6
F1 score	F1 score, also known as balanced F-score or F-measure, can be interpreted as a harmonic mean of precision and recall, where an F1 score reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for the F1 score is:  $F1 = 2 * TP / (2 * TP + FN + FP)$  Where "TP" is the number of true positives, "FN" is the number of false negatives, and "FP" is the number of false positives. F1 is by default calculated as 0.0 when there are no true positives, false negatives, nor false positives.	0.96045197740113

## Interpretations

Values close to 1 ( $>0.95$ ) were obtained for all the parameters, indicating that the Logistic Regression Model is a good classification model for the breast cancer dataset

## EXPERIMENT-2

**DATE: 23/01/24**

### Model Description and Input

The models used were:

- 1) Random Forests
- 2) Decision Trees
- 3) Support Vector Machine
- 4) XGBoost

The input was a Breast cancer dataset present as a default dataset in scikit learn (for practice purposes). The same steps as that of logistic regression were implemented in all the above 4 models and the various results obtained were compared.

### Results

PARAMETER	RANDOM FOREST	DECISION TREE	SVM	XGBOOST
Accuracy score (normalized)	0.972027972027972	0.8811188811188811	0.9370629370629371	0.972027972027972
Accuracy score (no normalization)	139	126	134	139
Area under the curve (AUC)	0.9969601677148846	0.8939203354297695	0.9842767295597485	0.9987421383647799
Sensitivity	0.9666666666666667	0.8444444444444444	0.9888888888888889	0.9666666666666667
Specificity	0.9811320754716981	0.9433962264150944	0.8490566037735849	0.9811320754716981

F1 score	0.9775280898 876404	0.89940828402 36687	0.9518716577 540107	0.977528089887 6404
Precision	0.9886363636 363636	0.96202531645 56962	0.9175257731 958762	0.988636363636 3636

## Interpretation

- It was observed that the runtime of the random forests code was greater compared to the others
- Even xgboost had a slightly greater runtime relative to decision trees and svm
- Decision trees gave a lower accuracy, AUC, Sensitivity and F1 score relative to the other models indicating that this model is not that suitable for predicting breast cancer
- Random Forests and XGBoost models gave exactly same values for all the metric (accuracy, precision, sensitivity etc) when the data was split into train and test sets using random state = 0 , changing the value of random state (say 40 for both models) gave different values.

## EXPERIMENT-3

**DATE: 24/01/24**

### Model Description and Input

The input was a Breast cancer dataset present as a default dataset in scikit learn (for practice purposes). The logistic regression code of Experiment 1 was modified to scale data and carry out 10 fold cross validation on the same as part of the pipeline. The pipeline (Pipeline from scikit-learn) includes two steps: standardization (StandardScaler) and logistic regression (LogisticRegression). StandardScaler transform was applied to standardize the input variables. The pipeline was further evaluated using a 10 fold cross validation - RepeatedStratifiedKFold. Each fold was trained on the scaled data

`n_splits=10`: Specifies the number of folds in the cross-validation. In this case, it's set to 10, meaning the dataset will be split into 10 parts, and the model will be trained and evaluated 10 times, each time using a different fold as the test set and the remaining data as the training set.

`n_repeats=3`: Specifies the number of times the cross-validation process will be repeated. It repeats the entire cross-validation process three times with different random splits.

`random_state=1`: Provides a seed for the random number generator. Setting a seed ensures reproducibility, meaning if you run the code with the same seed, you should get the same random splits each time.

### Results

```
Accuracy of pipeline: 0.977 (0.018)
Accuracy Score - default: 0.958041958041958
Accuracy Score - no normalization: 137
AUC: 0.9914046121593292
Sensitivity: 0.9666666666666667
Specificity: 0.9433962264150944
F1 score: 0.9666666666666667
Precision: 0.9666666666666667
```

	precision	recall	f1-score	support
0	0.94	0.94	0.94	53
1	0.97	0.97	0.97	90
accuracy			0.96	143
macro avg	0.96	0.96	0.96	143
weighted avg	0.96	0.96	0.96	143

```
Accuracy Score - default: 0.951048951048951
Accuracy Score - no normalization: 136
AUC: 0.9937106918238994
Sensitivity: 0.9444444444444444
Specificity: 0.9622641509433962
F1 score: 0.96045197740113
Precision: 0.9770114942528736
```

	precision	recall	f1-score	support
0	0.91	0.96	0.94	53
1	0.98	0.94	0.96	90
accuracy			0.95	143
macro avg	0.94	0.95	0.95	143
weighted avg	0.95	0.95	0.95	143



LHS -> Metrics obtained after standardization and cross validation

RHS -> Metrics obtained after experiment 1

## Interpretation

The pipeline accuracy was close to 1 indicating that the standardization done was proper. Slight increase in accuracy and sensitivity was observed. However, overall the model was not quite affected by the standardization, indicated by the similar values of the metrics pre and post standardization

## EXPERIMENT-4

**DATE: 25/01/24**

### Model Description and Input

Implemented 10 fold cross validation(without using the default function in sci-kit learn) and standardization on all the 5 classifiers and printed metrics for each fold, finally printing the mean and standard deviation of each of the 6 metrics. Input used was the breast cancer dataset present by default in sci-kit learn.

### Results

#### Mean and Standard deviation of Metrics over 10 fold cross validation for the 5 classifiers

##### LOGISTIC REGRESSION

+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+		
Accuracy	0.865972	0.274663
+-----+-----+-----+		
AUC	0.902248	0.285399
+-----+-----+-----+		
Sensitivity	0.883766	0.280223
+-----+-----+-----+		
Specificity	0.836678	0.271012
+-----+-----+-----+		
F1 score	0.875243	0.277244
+-----+-----+-----+		
Precision	0.867883	0.276302



+-----+-----+-----+

### **RANDOM FORESTS**

+-----+-----+-----+

| Metric | Mean | Standard Deviation |

+=====+=====+=====+

| Accuracy | 0.870785 | 0.276297 |

+-----+-----+-----+

| AUC | 0.898192 | 0.284401 |

+-----+-----+-----+

| Sensitivity | 0.883622 | 0.280178 |

+-----+-----+-----+

| Specificity | 0.849075 | 0.276065 |

+-----+-----+-----+

| F1 score | 0.878971 | 0.278497 |

+-----+-----+-----+

| Precision | 0.875487 | 0.279092 |

+-----+-----+-----+

### **DECISION TREES**

+-----+-----+-----+

| Metric | Mean | Standard Deviation |

+=====+=====+=====+

| Accuracy | 0.821172 | 0.261122 |

+-----+-----+-----+

| AUC | 0.817454 | 0.260927 |

+-----+-----+-----+

| Sensitivity | 0.832468 | 0.266304 |

+-----+-----+-----+

| Specificity | 0.80244 | 0.267702 |

+-----+-----+-----+

F1 score	0.838499	0.266066	
----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Precision	0.847754	0.272047	
-----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

### **SUPPORT VECTOR MACHINE**

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Metric		Mean		Standard Deviation	
--------	--	------	--	--------------------	--

+	=====	+	=====	+	=====	+
---	-------	---	-------	---	-------	---

Accuracy	0.830799	0.264151	
----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

AUC	0.887194	0.280768	
-----	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Sensitivity	0.888817	0.281768	
-------------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Specificity	0.734357	0.243995	
-------------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

F1 score	0.849673	0.269425	
----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Precision	0.815038	0.260516	
-----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

### **XGBOOST**

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Metric		Mean		Standard Deviation	
--------	--	------	--	--------------------	--

+	=====	+	=====	+	=====	+
---	-------	---	-------	---	-------	---

Accuracy	0.881949	0.279702	
----------	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

AUC	0.903572	0.285805	
-----	----------	----------	--

+	-----	+	-----	+	-----	+
---	-------	---	-------	---	-------	---

Sensitivity	0.888672	0.281725	
-------------	----------	----------	--

+-----+-----+-----+			
Specificity	0.870523	0.277784	
+-----+-----+-----+			
F1 score	0.887495	0.281156	
+-----+-----+-----+			
Precision	0.886626	0.28121	
+-----+-----+-----+			

## Interpretation

Each of the models gave similar results for the metrics, which might have been due to the near ideal nature of the dataset

## EXPERIMENT-5

**DATE: 30/01/24**

### Description and Input

This experiment involved analysis of the different datasets in SPARK (.csv files) in R for detection of the number of individuals affected by ASD and the number of unaffected individuals.

### Results

File	Total Individuals	ASD Affected	Non-ASD
core_descriptive_variables-2023-07-21.csv	132368	132368	0
basic_medical_screening-2023-07-21.csv	192557	79631	112926
cbcl_1_5-2023-07-21.csv	2835	2835	0
cbcl_6_18-2023-07-21.csv	10290	10290	0
dcdq-2023-07-21.csv	35097	35097	0
individuals_registration-2023-07-21.csv	328974	132139	196835

predicted_iq_experimental-2023-07-21.csv	132368	132368	0
rbsr-2023-07-21.csv	46517	46517	0
roles_2023-07-17.csv	132368	132368	0
scq-2023-07-21.csv	89678	63491	26187
srs-2_adult_self-2023-07-21.csv	1798	0	1798
vineland-3-2023-07-21.csv	23160	23159	1
background_history_sibling-2023-07-21.csv	19700	21	19679

Filename: asd\_noasd\_num\_basicmedical.r and asd\_noasd\_num\_multipliedatasets.r

## Interpretation

Only 3 datasets : individual registration, basic medical screening and scq has relevant number of both ASD affected and unaffected individuals, indicating that these datasets can probably used for future experiments

## Recommendations for Improvement

- 1) The above 3 mentioned datasets can be used as primary datasets for building model
- 2) For the subject IDs mentioned in these 3 datasets, relevant information (columns) can be extracted from the other datasets and added to a single data frame
- 3) Datasets can be merged on the basis of unique subject ID

# EXPERIMENT-6

**DATE: 01/02/24**

## Description and Input

This experiment involved analysis of the basic medical screening dataset in SPARK (.csv file) in R for detection of the number of affected and unaffected individuals for the following disorders:

- 1) ASD
- 2) ADHD
- 3) ODD
- 4) OCD
- 5) Schizophrenia

## Results

**TOTAL NUMBER OF INDIVIDUALS -> 192557** , Filename : asd\_comorbidities.r

Disease	Non-Affected Individuals	Affected Individuals
ASD	112926	79631
ADHD	148304	44253
ODD	185436	7121
OCD	179556	13001
Schizophrenia	190794	1763

## EXPERIMENT-7

**DATE: 02/02/24**

### Description and Input

This experiment involved analysis of the basic medical screening dataset in SPARK (.csv file) in R for detection of the following conditions pertaining ASD and ADHD:

- 1) Individuals affected by ASD only and no other comorbidities
- 2) Individuals affected by ADHD only and no other comorbidities
- 3) Individuals affected by both ASD and ADHD
- 4) Individuals unaffected by any of the 2 disorders

### Results

**TOTAL NUMBER OF INDIVIDUALS -> 192557**

ASD only	ADHD only	Both ASD and ADHD	Unaffected by any of the 2 disorders
49804	14426	29827	98500

Filename: asd\_adhd\_numbers.r

# EXPERIMENT-8

**DATE: 03/02/24**

## Description and Input

This experiment involved analysis of the basic medical screening dataset in SPARK (.csv file) and filtering out the columns that will be used to develop the ML model (based on certain criteria inferred from literature).

## Results

**TOTAL NUMBER OF INDIVIDUALS -> 192557**

### IGNORED COLUMNS:

- Ignored the columns containing various SPARK ids because the dataset will be used as a whole for model building and not use individual ids
  - subject\_sp\_id
  - respondent\_sp\_id
  - family\_sf\_id
  - biomother\_sp\_id
  - biofather\_sp\_id
  - current\_depend\_adult
- Ignored the columns concerning other comorbidities:
  - behav\_conduct
  - behav\_intermitt\_explos
  - behav\_odd
- Current\_depend\_adult
- Flag
- The following columns are broad columns having sub-columns ( branching qs ) :- They are probably placing 1/ null on the basis of answers to all the sub questions but it is vague what is the basis of choosing the final value
  - Attn\_behav
  - Birth\_def\_cns -> Branching Question: Brain and spinal cord birth defects
  - Birth\_def\_bone
  - birth\_def\_fac
  - Birth\_def\_gastro



- Birth\_def\_thorac (heart or lung)
  - Birth\_def\_urogen
  - Dev\_lang
  - Gen\_test
  - med\_cond\_birth
  - med\_cond\_birth\_def
  - med\_cond\_growth
  - med\_cond\_neuro
  - Med\_cond\_visaud
  - mood OCD
- Birth defects: (defects not impacting asd/adhd)
  - Birth\_def\_bone\_club
  - birth\_def\_bone\_miss
  - birth\_def\_bone\_polydact
  - birth\_def\_bone\_spine
  - birth\_def\_cleft\_lip
  - Birth\_def\_cleft\_palate
  - Birth\_def\_cns\_myelo
  - birth\_def\_gi\_esoph\_atres
  - birth\_def\_gi\_hirschprung
  - birth\_def\_gi\_intest\_malrot
  - Birth\_def\_gi\_pylor\_sten
  - birth\_def\_thorac\_heart
  - Birth\_def\_thorac\_lung
  - birth\_def\_urogen\_hypospad
  - birth\_def\_urogen\_renal
  - birth\_def\_urogen\_renal\_agen
  - birth\_def\_urogen\_uter\_agen
  - Birth\_def\_oth\_calc
  - birth\_ivh
  - Birth\_oth\_calc
- Eating\_probs - Not professionally diagnosed so probably not a definitive measure
- Eating\_disorder - Not directly associated with ASD
- Etoh\_subst
- gen\_dx\_oth\_calc\_self\_report
- Genetic testing:
  - gen\_test\_cgh\_cma
  - gen\_test\_chrom\_karyo

- 
- gen\_test\_ep
  - gen\_test\_fish\_angel
  - gen\_test\_fish\_digeorge
  - gen\_test\_fish\_williams
  - gen\_test\_fish\_oth
  - gen\_test\_frax
  - gen\_test\_id
  - gen\_test\_mecp2
  - gen\_test\_nf1
  - gen\_test\_noonan
  - gen\_test\_pten
  - gen\_test\_tsc
  - gen\_test\_unknown
  - gen\_test\_wes
  - gen\_test\_wgs
  - Gen\_test\_oth\_calc
  - Calculated variable: individual has entered in text field an "Other" growth condition not represented in current coding
  - growth\_low\_wt
  - growth\_macroceph
  - growth\_microceph
  - growth\_obes
  - growth\_short
  - Growth\_oth\_calc
  - prev\_study\_calc
  - Eval\_year
  - neuro\_inf
  - Neuro\_lead
  - neuro\_sz
  - neuro\_tbi

- Neuro\_oth\_calc
- pers\_dis
- prev\_study\_oth\_calc
- psych\_oth\_calc
- Schiz
- visaud\_blind
- visaud\_catar
- visaud\_deaf
- Visaud\_strab
- sleep\_dx
- sleep\_eat\_toilet
- sleep\_probs
- tics

## CHOSEN COLUMNS:

- 1) Sex -> To know which gender is more susceptible to a particular disorder
- 2) Asd -> Professional diagnosis results
- 3) Age\_at\_eval\_months -> To know the age distribution of each disorder
- 4) Age\_at\_eval\_years -> To know the age distribution of each disorder ( we can use any one of 3 or 4 -
- 5) Behav\_adhd -> professional diagnosis of adhd
- 6) Birth defects:
  - a) Birth\_def\_cns\_brain -> Brain malformation/abnormality (shown on MRI)
  - b) Birth\_def\_thorac\_cdh -> Congenital diaphragmatic hernia - Increased Risk:  
Studies have shown that children with CDH have a higher risk of developing Autism Spectrum Disorder (ASD) compared to the general population. The estimated prevalence of ASD in CDH survivors is around 9-14%, compared to 1.5% in the general population. Although there is research on the increased risk of ASD in CDH, there is currently no established direct link between CDH and Attention Deficit Hyperactivity Disorder (ADHD).
  - c) Birth\_etoh\_subst -> Fetal Alcohol Syndrome, alcohol or drug exposure in mother's pregnancy

- d) Birth\_oxygen - > Insufficient oxygen at birth with NICU stay - Studies suggest a higher risk of ASD and ADHD in children who experienced HIE compared to the general population.
- e) Birth\_pg\_inf -> Serious prenatal infection (for example, German measles) - Studies suggest a modest increase in the risk of ASD in children whose mothers contracted certain severe infections during pregnancy
- f) Birth\_prem -> Premature birth (delivery before 37 weeks) - Premature birth, delivering before 37 weeks of pregnancy, is indeed linked to an increased risk of Autism Spectrum Disorder (ASD) and Attention Deficit Hyperactivity Disorder (ADHD). The earlier the birth, the higher the risk, according to multiple studies
- 7) cog\_med- Cognitive delays or impairment due to another medical condition or exposure (For example, brain injury, stroke, lead poisoning, FAS, HIV, radiation, hydrocephalus, brain tumor, drug effects, etc.) —> Too many factors so might make it confusing
- 8) Developmental properties:
  - a) Dev\_id - Intellectual disability, cognitive impairment, global developmental delay, or borderline intellectual functioning
  - b) Dev\_lang\_dis - Language difficulties are core features of ASD for many individuals. While not a core symptom of ADHD, individuals with ADHD can sometimes experience language difficulties.
  - c) Dev\_ld - Learning disability (LD, learning disorder, including reading, written expression, math, or NVLD (Nonverbal learning disability))
  - d) Dev\_motor - Motor delay (e.g., delay in walking) or developmental coordination disorder - Both motor delays and DCD can co-occur with ASD and ADHD.
  - e) Dev\_mutism - Mutism, particularly SM, can co-occur with ASD and ADHD.
  - f) Dev\_soc\_prag - Social (Pragmatic) Communication Disorder - SCD can co-occur with both ASD and ADHD.
  - g) Dev\_speech - Speech articulation problems are more common in individuals with ASD compared to the general population. Not much in ADHD
- 9) Encopres - Encopresis (has bowel accidents beyond age expected). Studies suggest children with ASD and ADHD might be at higher risk for encopresis compared to the general population.
- 10) Enures - Enuresis (wets self beyond age expected). Children with ASD are 2-3 times more likely to experience enuresis compared to typically developing children

11) Feeding\_dx - Feeding/eating problems - Children with ASD are five times more likely to experience feeding difficulties compared to typically developing peers.

12) Parent-reported genetic testing:

a) Gen\_test\_aut\_dd - Parent reported Genetic testing for autism or developmental delay genes

13) Gest\_age - How many weeks (gestational age) was a child/dependent when he/she was born? - Premature birth (born before 37 weeks) is associated with a slightly increased risk of developing both ASD and ADHD.

14) Mood disorders: - usually associated with ASD / ADHD

a) Mood\_anx - Anxiety disorder, such as panic, phobia, agoraphobia, or generalized anxiety disorder (GAD) except for social anxiety

b) Mood\_bipol - Bipolar (Manic-Depressive) Disorder

c) Mood\_dep - Depression or dysthymia

d) Mood\_dmd - Disruptive Mood Dysregulation Disorder

e) Mood\_hoard - Hoarding

f) Mood\_or\_anx - Obsessive-Compulsive Disorder

g) Mood\_sep\_anx - Separation Anxiety

h) Mood\_soc\_anx - Social Anxiety Disorder/Social Phobia

prev\_study\_agre

prev\_study\_asc

prev\_study\_charge

prev\_study\_earli

prev\_study\_marbles

prev\_study\_mssng

prev\_study\_seed

prev\_study\_ssc

prev\_study\_vip

# EXPERIMENT-9

**DATE: 17/02/24**

## Description and Input

We developed a XGBoost Classifier on the basic medical screening dataset to make a 3-way classification model that classifies patients into individuals with ASD only, ADHD only or both.

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** ['sex']

**Numerical columns :** ['age\_at\_eval\_months', 'age\_at\_eval\_years', 'birth\_def\_cns\_brain', 'birth\_def\_thorac\_cdh', 'birth\_etoh\_subst', 'birth\_oxygen', 'birth\_pg\_inf', 'birth\_prem', 'cog\_med', 'dev\_id', 'dev\_lang\_dis', 'dev\_ld', 'dev\_motor', 'dev\_mutism', 'dev\_soc\_prag', 'dev\_speech', 'eating\_probs', 'eating\_disorder', 'encopres', 'enures', 'feeding\_dx', 'gen\_test\_aut\_dd', 'gest\_age', 'mood\_anx', 'mood\_bipol', 'mood\_dep', 'mood\_dmd', 'mood\_hoard', 'mood\_or\_anx', 'mood\_sep\_anx', 'mood\_soc\_anx', 'prev\_study\_agre', 'prev\_study\_asc', 'prev\_study\_charge', 'prev\_study\_earli', 'prev\_study\_marbles', 'prev\_study\_mssng', 'prev\_study\_seed', 'prev\_study\_ssc', 'prev\_study\_vip', 'sleep\_dx', 'sleep\_eat\_toilet', 'sleep\_probs']

## Results

Df size after removing unaffected people = 94057

Number of X features = 44

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

```
+-----+-----+-----+
| Metric   | Mean | Standard Deviation |
+=====+=====+=====+
```

Accuracy	0.686094	0.00584678	
+-----+-----+-----+			
Sensitivity	0.671497	0.0134625	
+-----+-----+-----+			
Specificity	0.939345	0.00220678	
+-----+-----+-----+			
F1 score	0.682309	0.00617388	
+-----+-----+-----+			
Precision	0.681323	0.00621106	
+-----+-----+-----+			

Filename: model\_3way\_basicmedical\_1.py

## Interpretation

The classification of the numerical and categorical columns was incorrect as all the columns with Boolean values became a part of the Numerical columns while they should have been a part of the categorical columns. The ROC-AUC value is not being calculated for a 3 way classification model.

## Recommendations for Improvement

- The classification of the numerical and categorical columns should be improved - by manually observing the numerical and categorical columns
- Modifying the chosen features - some features are important but not part of the model

# EXPERIMENT-10

**DATE: 20/02/24**

## Description and Input

We modified the above XGBoost Classifier on the basic medical screening dataset to make a 3-way classification model that classifies patients into individuals with ASD only, ADHD only or both by making changes in the features being used, also modifying the columns (categorical and numerical) used by the column transformer

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** ['birth\_def\_bone\_club', 'birth\_def\_bone\_miss', 'birth\_def\_bone\_polydact', 'birth\_def\_bone\_spine', 'birth\_def\_cleft\_lip', 'birth\_def\_cleft\_palate', 'birth\_def\_cns\_brain', 'birth\_def\_cns\_myelo', 'birth\_def\_gi\_esoph\_atres', 'birth\_def\_gi\_hirschprung', 'birth\_def\_gi\_intest\_malrot', 'birth\_def\_gi\_pylor\_sten', 'birth\_def\_thorac\_cdh', 'birth\_def\_thorac\_heart', 'birth\_def\_thorac\_lung', 'birth\_def\_urogen\_hypospad', 'birth\_def\_urogen\_renal', 'birth\_def\_urogen\_renal\_agen', 'birth\_def\_urogen\_uter\_agen', 'birth\_etoh\_subst', 'birth\_ivh', 'birth\_oxygen', 'birth\_pg\_inf', 'birth\_prem', 'cog\_med', 'dev\_id', 'dev\_lang\_dis', 'dev\_ld', 'dev\_motor', 'dev\_mutism', 'dev\_soc\_prag', 'dev\_speech', 'eating\_probs', 'eating\_disorder', 'encopres', 'enures', 'feeding\_dx', 'growth\_low\_wt', 'growth\_macroceph', 'growth\_microceph', 'growth\_obes', 'growth\_short', 'mood\_anx', 'mood\_bipol', 'mood\_dep', 'mood\_dmd', 'mood\_hoard', 'mood\_or\_anx', 'mood\_sep\_anx', 'mood\_soc\_anx', 'neuro\_inf', 'neuro\_lead', 'neuro\_sz', 'neuro\_tbi', 'sleep\_dx', 'sleep\_eat\_toilet', 'sleep\_probs', 'visaud\_blind', 'visaud\_catar', 'visaud\_deaf', 'visaud\_strab']

**Numerical columns :** ['gest\_age']

## Results

Df size after removing unaffected people = 94057

Number of X features = 62

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS



## MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

+-----+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+=====+		
Accuracy	0.650914	0.00487393
+-----+-----+-----+-----+		
Sensitivity	0.538473	0.00984392
+-----+-----+-----+-----+		
Specificity	0.928407	0.00376441
+-----+-----+-----+-----+		
F1 score	0.641011	0.00519251
+-----+-----+-----+-----+		
Precision	0.643825	0.00509689
+-----+-----+-----+-----+		

Filename: model\_3way\_basicmedical\_2.py

## Interpretation

The low accuracy but high sensitivity might indicate that the features being used are good for predicting ASD only → ADHD is not being predicted well.

## Recommendations for Improvement

- 1) Reading up about the features being used and finding out what can be eliminated - features should be decreased
- 2) Try ordinal encoding instead of one hot on all the categorical columns
- 3) Try recursive feature elimination to work on features
- 4) Gestational age column - replace NAN with 40
- 5) EDA on categorical columns - as a bar plot
- 6) To make it inclusive of normal individuals - include bg datasets - bg child, sibling etc.

# EXPERIMENT-11

**DATE: 27/02/24**

## Description and Input

I carried out Exploratory data analysis on the chosen features to understand which of the features are truly relevant. I plotted a histogram for the numerical column (gestational age) and created a table for the categorical ones.

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** ['birth\_def\_bone\_club', 'birth\_def\_bone\_miss', 'birth\_def\_bone\_polydact', 'birth\_def\_bone\_spine', 'birth\_def\_cleft\_lip', 'birth\_def\_cleft\_palate', 'birth\_def\_cns\_brain', 'birth\_def\_cns\_myelo', 'birth\_def\_gi\_esoph\_atres', 'birth\_def\_gi\_hirschprung', 'birth\_def\_gi\_intest\_malrot', 'birth\_def\_gi\_pylor\_sten', 'birth\_def\_thorac\_cdh', 'birth\_def\_thorac\_heart', 'birth\_def\_thorac\_lung', 'birth\_def\_urogen\_hypospad', 'birth\_def\_urogen\_renal', 'birth\_def\_urogen\_renal\_agen', 'birth\_def\_urogen\_uter\_agen', 'birth\_etoh\_subst', 'birth\_ivh', 'birth\_oxygen', 'birth\_pg\_inf', 'birth\_prem', 'cog\_med', 'dev\_id', 'dev\_lang\_dis', 'dev\_ld', 'dev\_motor', 'dev\_mutism', 'dev\_soc\_prag', 'dev\_speech', 'eating\_probs', 'eating\_disorder', 'encopres', 'enures', 'feeding\_dx', 'growth\_low\_wt', 'growth\_macroceph', 'growth\_microceph', 'growth\_obes', 'growth\_short', 'mood\_anx', 'mood\_bipol', 'mood\_dep', 'mood\_dmd', 'mood\_hoard', 'mood\_or\_anx', 'mood\_sep\_anx', 'mood\_soc\_anx', 'neuro\_inf', 'neuro\_lead', 'neuro\_sz', 'neuro\_tbi', 'sleep\_dx', 'sleep\_eat\_toilet', 'sleep\_probs', 'visaud\_blind', 'visaud\_catar', 'visaud\_deaf', 'visaud\_strab']

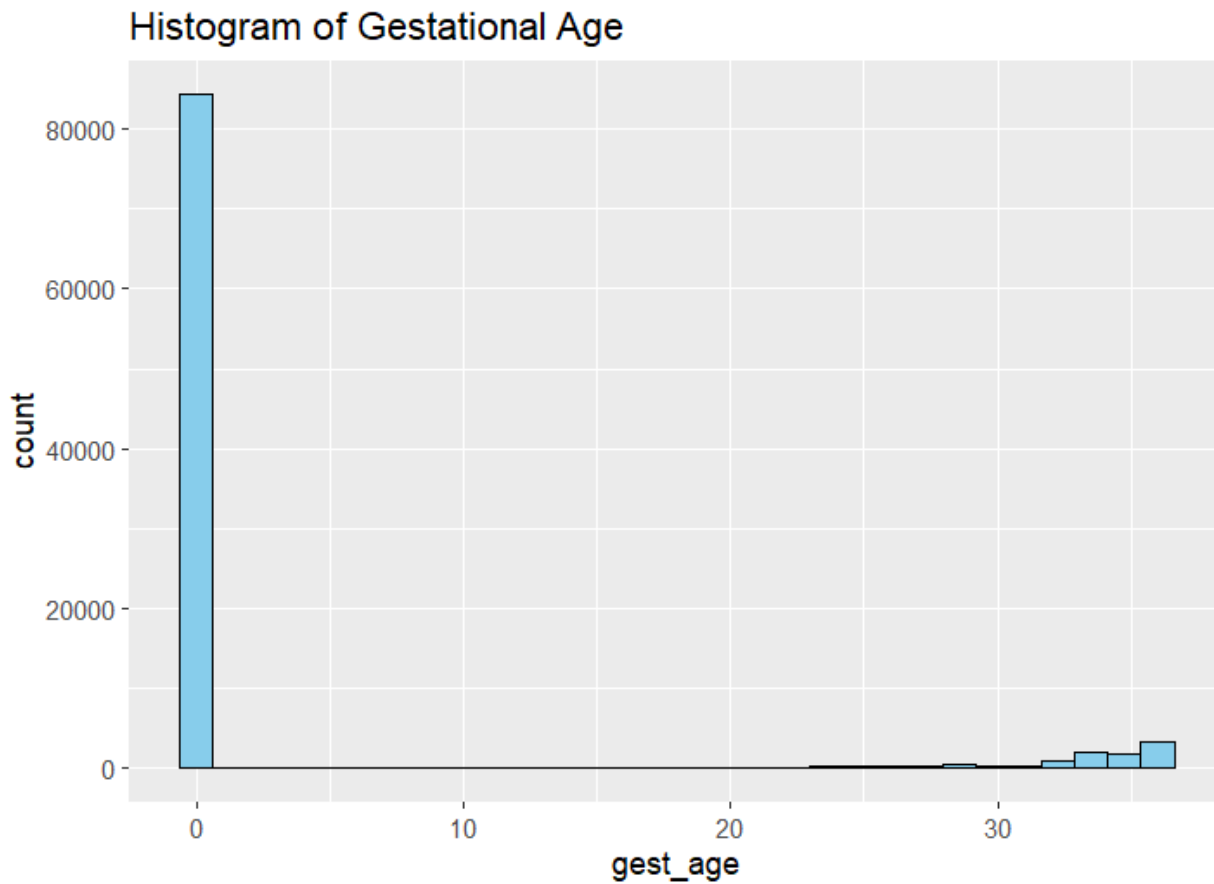
**Numerical columns :** ['gest\_age']

## Results

Df size after removing unaffected people = 94057


Number of X features = 62

### Results for gestational age:




### Results for categorical columns


Column_Name	Value_1_Count	Value_0_Count
birth_def_bone_club	200	93857
birth_def_bone_miss	356	93701
birth_def_bone_polydact	73	93984
birth_def_bone_spine	154	93903



birth_def_cleft_lip	106	93951
birth_def_cleft_palate	147	93910
birth_def_cns_brain	280	93777
birth_def_cns_myelo	63	93994
birth_def_gi_esoph_atres	32	94025
birth_def_gi_hirschprung	15	94042
birth_def_gi_intest_malrot	45	94012
birth_def_gi_pylor_sten	38	94019
birth_def_thorac_cdh	23	94034
birth_def_thorac_heart	664	93393
birth_def_thorac_lung	58	93999
birth_def_urogen_hypospa d	257	93800
birth_def_urogen_renal	83	93974



birth_def_urogen_renal_age	20	94037
birth_def_urogen_uter_age	5	94052
birth_etoh_subst	1092	92965
birth_ivh	637	93420
birth_oxygen	5461	88596
birth_pg_inf	228	93829
birth_prem	10380	83677
cog_med	4348	89709
dev_id	16516	77541
dev_lang_dis	41394	52663
dev_ld	19937	74120
dev_motor	13796	80261
dev_mutism	1013	93044



dev_soc_prag	13506	80551
dev_speech	21610	72447
eating_probs	29433	64624
eating_disorder	1840	92217
encopres	5850	88207
enures	6368	87689
feeding_dx	11593	82464
growth_low_wt	2974	91083
growth_macroceph	2473	91584
growth_microceph	796	93261
growth_obes	2302	91755
growth_short	2015	92042
mood_anx	25518	68539

mood_bipol	4814	89243
mood_dep	16045	78012
mood_dmd	2612	91445
mood_hoard	1770	92287

Filename: eda\_model.r

## Interpretation

Many columns have value <1000 for Value=1 and can be ignored. These are:

'birth\_def\_bone\_club', 'birth\_def\_bone\_miss', 'birth\_def\_bone\_polydact',  
 'birth\_def\_bone\_spine', 'birth\_def\_cleft\_lip', 'birth\_def\_cleft\_palate',  
 'birth\_def\_cns\_brain', 'birth\_def\_cns\_myelo', 'birth\_def\_gi\_esoph\_atres',  
 'birth\_def\_gi\_hirschprung', 'birth\_def\_gi\_intest\_malrot', 'birth\_def\_gi\_pylor\_sten',  
 'birth\_def\_thorac\_cdh', 'birth\_def\_thorac\_heart', 'birth\_def\_thorac\_lung',  
 'birth\_def\_urogen\_hypospad', 'birth\_def\_urogen\_renal', 'birth\_def\_urogen\_renal\_agen',  
 'birth\_def\_urogen\_uter\_agen', 'birth\_ivh', 'birth\_pg\_inf', 'growth\_microceph', 'neuro\_inf',  
 'neuro\_lead', 'neuro\_tbi', 'visaud\_blind', 'visaud\_catar'

## **EXTRA - FEATURES TO BE USED**

sex

birth\_def\_bone\_club

birth\_def\_bone\_miss

birth\_def\_bone\_polydact

birth\_def\_bone\_spine

birth\_def\_cleft\_lip

birth\_def\_cleft\_palate

birth\_def\_cns\_brain

birth\_def\_cns\_myelo

birth\_def\_fac

birth\_def\_gi\_esoph\_atres

birth\_def\_gi\_hirschprung

birth\_def\_gi\_intest\_malrot

birth\_def\_gi\_pylor\_sten

birth\_def\_thorac\_cdh

birth\_def\_thorac\_heart

birth\_def\_thorac\_lung

birth\_def\_urogen\_hypospad


birth\_def\_urogen\_renal

birth\_def\_urogen\_renal\_agen


birth\_def\_urogen\_uter\_agen

birth\_etoh\_subst





birth\_ivh  
birth\_oth\_calc  
birth\_oxygen  
birth\_pg\_inf  
cog\_med  
dev\_id  
dev\_lang  
dev\_lang\_dis  
dev\_ld  
dev\_motor  
dev\_mutism  
dev\_soc\_prag  
dev\_speech  
eating\_probs  
eating\_disorder  
  
feeding\_dx  
  
gest\_age  
growth\_low\_wt  
growth\_macroceph  
growth\_microceph  
growth\_obes  
growth\_short  
  
med\_cond\_birth  
med\_cond\_birth\_def  
med\_cond\_growth



```
med_cond_neuro
med_cond_visaud
```

```
neuro_inf
neuro_lead
neuro_sz
neuro_tbi
```

```
features = [
    "sex", "gest_age", "eating_probs", "feeding_dx",
    "med_cond_birth", "birth_oth_calc",
    "med_cond_birth_def", #"birth_def_oth_calc" is not used because all are nans
    "med_cond_growth", "growth_oth_calc",
    "med_cond_neuro", "med_cond_visaud"
]

# use this set of features for SPARK experiments
features = [
    "sex", 'mother_highest_education', 'father_highest_education',
    'annual_household_income',
    "smiled_age_mos", "sat_wo_support_age_mos", "crawled_age_mos", "walked_age_mos",
    "fed_self_spoon_age_mos", "used_words_age_mos", "combined_words_age_mos",
    "combined_phrases_age_mos", "bladder_trained_age_mos", "bowel_trained_age_mos",
    "hand",
    "twin_mult_birth", 'num_asd_parents', 'num_asd_siblings'
]
```

# EXPERIMENT-12

**DATE: 28/02/24**

## Description and Input

We modified the above XGBoost Classifier on the basic medical screening dataset to make a 3-way classification model that classifies patients into individuals with ASD only, ADHD only or both by replacing nan values in the gest\_age column with 40. Some features are removed as well.

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** ['birth\_etoh\_subst', 'birth\_oxygen', 'birth\_prem', 'cog\_med', 'dev\_id', 'dev\_lang\_dis', 'dev\_ld', 'dev\_motor', 'dev\_mutism', 'dev\_soc\_prag', 'dev\_speech', 'eating\_probs', 'eating\_disorder', 'encopres', 'enures', 'feeding\_dx', 'growth\_low\_wt', 'growth\_macroceph', 'growth\_obes', 'growth\_short', 'mood\_anx', 'mood\_bipol', 'mood\_dep', 'mood\_dmd', 'mood\_hoard', 'mood\_or\_anx', 'mood\_sep\_anx', 'mood\_soc\_anx', 'neuro\_sz', 'sleep\_dx', 'sleep\_eat\_toilet', 'sleep\_probs', 'visaud\_deaf', 'visaud\_strab']

**Numerical columns :** ['gest\_age']

## Results

Number of X features = 35

Df size after removing unaffected people = 94057

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

+-----+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+=====+		
Accuracy	0.649957	0.00502995
+-----+-----+-----+-----+		
Sensitivity	0.53168	0.0119204
+-----+-----+-----+-----+		



Specificity	0.928131	0.00273799
-------------	----------	------------

+-----+	+-----+	+-----+
---------	---------	---------

F1 score	0.639835	0.0055437
----------	----------	-----------

+-----+	+-----+	+-----+
---------	---------	---------

Precision	0.642732	0.00539781
-----------	----------	------------

+-----+	+-----+	+-----+
---------	---------	---------

Filename: model\_3way\_basicmedical\_3.py

# EXPERIMENT-13

**DATE: 29/02/24**

## Description and Input

We modified the above XGBoost Classifier on the basic medical screening dataset to make a 3-way classification model that classifies patients into individuals with ASD only, ADHD only or both by using 11 features.

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** ['birth\_etoh\_subst', 'birth\_oxygen', 'birth\_prem', 'cog\_med', 'dev\_id', 'dev\_lang\_dis', 'dev\_ld', 'dev\_motor', 'dev\_mutism', 'dev\_soc\_prag', 'dev\_speech', 'eating\_probs', 'eating\_disorder', 'encopres', 'enures', 'feeding\_dx', 'growth\_low\_wt', 'growth\_macroceph', 'growth\_obes', 'growth\_short', 'mood\_anx', 'mood\_bipol', 'mood\_dep', 'mood\_dmd', 'mood\_hoard', 'mood\_or\_anx', 'mood\_sep\_anx', 'mood\_soc\_anx', 'neuro\_sz', 'sleep\_dx', 'sleep\_eat\_toilet', 'sleep\_probs', 'visaud\_deaf', 'visaud\_strab']

**Numerical columns :** ['gest\_age']


## Results

Number of X features = 11

Df size after removing unaffected people = 94057

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

Metric	Mean	Standard Deviation
Accuracy	0.563212	0.00516725
Sensitivity	0.559269	0.0169223
Specificity	0.906569	0.00333269



```
+-----+-----+-----+
| F1 score | 0.494427 |    0.00511255 |
+-----+-----+-----+
| Precision | 0.537941 |    0.00903289 |
+-----+-----+-----+
```

Filename: model\_3way\_basicmedical\_4.py

## EXPERIMENT-14

**DATE: 04/03/24**

### Description and Input

We modified the above XGBoost Classifier on the basic medical screening dataset to make a 2-way classification model that classifies patients into individuals with ADHD or not using 11 features. We also tried out the same code using other classifiers used earlier to compare the values of the different metrics. We tried this with and without class balance.

#### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** 'sex', 'eating\_probs', 'feeding\_dx',

'med\_cond\_birth', 'birth\_oth\_calc',

'med\_cond\_birth\_def',

'med\_cond\_growth', 'growth\_oth\_calc',

'med\_cond\_neuro', 'med\_cond\_visaud'

**Numerical columns :** ['gest\_age']

### Results

Number of X features = 11

Number of samples = 192557

Number of individuals affected by ADHD: 44253

Number of unaffected individuals: 148304

#### CASE 1: WITHOUT CLASS BALANCE

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD

Accuracy	0.7698 3459	0.000 74839 3	0.769 367	0.000 65391 3	0.7687 85	0.000 8960 79	0.7688 11	0.000 80017 2
Sensitivity	0.9836 21	0.001 22201	0.982 819	0.001 05313	0.9817 6	0.001 4927 2	0.9800 95	0.000 96415 4
Specificity	0.0533 75	0.004 34085	0.054 0304	0.002 16125	0.0550 472	0.002 859	0.0607 419	0.002 63818
F1 score	0.6907 39	0.001 75082	0.690 699	0.000 86351 9	0.6907 06	0.001 0792 5	0.6926 63	0.001 20247
Precision	0.7116 27	0.003 81884	0.709 654	0.002 96767	0.7073 54	0.003 8157 8	0.7084 62	0.003 35527
ROC-AUC Score	0.6507 92	0.005 35774	0.648 02	0.004 88875	0.6425 75	0.006 2481 3	0.6448 86	0.005 97992

Filename: model\_2way\_basicmedical\_adhd\_1.py

### **CASE 2: WITH CLASS BALANCE**

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.6185 57	0.005 2565	0.614 908	0.006 80877	0.6141 84	0.004 7184 1	0.6158 68	0.004 74355
Sensitivity	0.7231 6	0.005 38263	0.762 389	0.006 75764	0.7301 65	0.005 9262	0.7275 21	0.006 19465
Specificity	0.5139 54	0.009 87442	0.467 427	0.016 2587	0.4982 04	0.008 9733 3	0.5042 15	0.009 68887
F1 score	0.6143	0.005	0.606	0.007	0.6089	0.005	0.6109	0.005



	19	59345	282	91234	06	0261 7	98	08422
Precision	0.6239 9	0.005 17802	0.625 864	0.006 26792	0.6206 89	0.004 7498 2	0.6219 63	0.004 71774
ROC-AUC Score	0.6485 31	0.005 03725	0.646 848	0.005 28093	0.6434 74	0.005 0575 5	0.6444 31	0.004 93622

Filename: model\_2way\_basicmedical\_adhd\_2.py

# EXPERIMENT-15

**DATE: 05/03/24**

## Description and Input

We modified the above XGBoost Classifier on the basic medical screening dataset to make a 2-way classification model that classifies patients into individuals with ASD or not using 11 features. We also tried out the same code using other classifiers used earlier to compare the values of the different metrics. We tried this with and without class balance.

### COLUMNS CHOSEN FOR FEATURES ARE:

**Categorical columns :** 'sex', 'eating\_probs', 'feeding\_dx',

'med\_cond\_birth', 'birth\_oth\_calc',

'med\_cond\_birth\_def',

'med\_cond\_growth', 'growth\_oth\_calc',

'med\_cond\_neuro', 'med\_cond\_visaud'

**Numerical columns :** ['gest\_age']

## Results

Number of X features = 11

Number of samples = 192557

Number of individuals affected by ASD 79631

Number of unaffected individuals: 112926

### CASE 1: WITHOUT CLASS BALANCE

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD

Accuracy	0.8172 33	0.001 86797	0.812 035	0.002 14957	0.8160 28	0.001 9088 5	0.8166 05	0.001 71887
Sensitivity	0.9165 65	0.001 84422	0.934 107	0.001 60649	0.9173	0.001 9256 5	0.9169 01	0.002 01713
Specificity	0.6763 7	0.004 60801	0.638 922	0.006 08115	0.6724 14	0.004 5736 2	0.6743 73	0.004 39409
F1 score	0.8129 45	0.002 00906	0.805 614	0.002 44153	0.8115 67	0.002 0475 7	0.8122 31	0.001 84546
Precision	0.8215 25	0.001 79507	0.821 631	0.001 7688	0.8206 16	0.001 8593 6	0.8210 41	0.001 68447
ROC-AUC Score	0.8655 2	0.001 82675	0.861 471	0.001 93155	0.8634 16	0.002 0814 5	0.8645	0.001 87666

Filename: model\_2way\_basicmedical\_asd\_1.py

### **CASE 2: WITH CLASS BALANCE**

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.7974 53	0.002 90681	0.790 17	0.002 76503	0.7960 34	0.002 8530 5	0.7967	0.002 99884
Sensitivity	0.9104 12	0.002 06959	0.928 52	0.002 51392	0.9110 52	0.002 1763 1	0.9113 54	0.002 12653
Specificity	0.6844 95	0.004 55889	0.651 819	0.004 26426	0.6810 16	0.004 4453 5	0.6820 46	0.004 70103

F1 score	0.7948 34	0.003 01352	0.786 073	0.002 88724	0.7932 98	0.002 9579 6	0.7939 9	0.003 1119
Precision	0.8134 54	0.002 67131	0.814 233	0.002 68003	0.8125 77	0.002 6507 9	0.8131 69	0.002 75482
ROC-AUC Score	0.8656 77	0.001 96406	0.861 935	0.002 1349	0.8636 16	0.002 0229 5	0.8645 49	0.002 07078

Filename: model\_2way\_basicmedical\_asd\_2.py

# EXPERIMENT-16

**DATE: 06/03/24**

## Description and Input

Analysis of the features being used as model input from the different datasets and finding out which feature belongs to which dataset and which are common to multiple datasets

### DATASETS CHOSEN FOR MODEL:

- 1) df1 = basic\_medical\_screening-2023-07-21.csv
- 2) df2 = background\_history\_adult-2023-07-21.csv
- 3) df3 = background\_history\_child-2023-07-21.csv
- 4) df4 = background\_history\_sibling-2023-07-21.csv
- 5) df5 = roles\_2023-07-17.csv
- 6) df6 = individuals\_registration-2023-07-21.csv

## Results

SL.NO	FEATURE	DATASETS TO WHICH IT BELONGS
1)	sex	All 6
2)	mother_highest_education	df2, df3, df4
3)	father_highest_education	df2, df3, df4
4)	annual_household_income	df2, df3, df4
5)	smiled_age_mos	df3, df4
6)	sat_wo_support_age_mos	df3, df4

7)	crawled_age_mos	df3, df4
8)	fed_self_spoon_age_mos	df3, df4
9)	used_words_age_mos	df3, df4
10)	combined_words_age_mos	df3, df4
11)	combined_phrases_age_mos	df3, df4
12)	bladder_trained_age_mos	df3, df4
13)	bowel_trained_age_mos	df3, df4
14)	twin_mult_birth	df3, df4
15)	num_asd_parents	df6
16)	num_asd_siblings	df6

# EXPERIMENT-17

**DATE: 13/03/24**

## Description and Input

Successfully implemented the 2-way classification model for ASD using 4 datasets and XGBoost classifier - with and without Class balancing.

### DATASETS CHOSEN FOR MODEL:

- 1) df1 = basic\_medical\_screening-2023-07-21
- 2) df2 = background\_history\_child-2023-07-21
- 3) df3 = background\_history\_sibling-2023-07-21
- 4) df4 = individuals\_registration-2023-07-21

## Results

Number of X features = 28, Number of samples = 65897 (without class balance), 38360 (with class balance)

ASD- affected = 46717, Unaffected = 19180

METRIC	XGBOOST		XGBOOST BALANCED	
	MEAN	SD	MEAN	SD
Accuracy	0.836396	0.00511932	0.814964	0.00758401
Sensitivity	0.682586	0.0126251	0.830918	0.00887275
Specificity	0.899544	0.0035826	0.799009	0.00771655
F1 score	0.834504	0.00531861	0.814914	0.00758268
Precision	0.833499	0.00540914	0.8153	0.00761652
ROC-AUC Score	0.895154	0.00558414	0.891596	0.00732261

Filename: model\_2way\_asd\_multidata.py & model\_2way\_asd\_multidata\_balanced.py

# EXPERIMENT-18

**DATE: 13/03/24**

## Description and Input

Successfully implemented the 2-way classification model for ADHD using 4 datasets and XGBoost classifier - with and without Class balancing.

### DATASETS CHOSEN FOR MODEL:

1. df1 = basic\_medical\_screening-2023-07-21
2. df2 = background\_history\_child-2023-07-21
3. df3 = background\_history\_sibling-2023-07-21
4. df4 = individuals\_registration-2023-07-21

## Results

Number of X features = 28, Number of samples = 65897 (without class balance), 40188 (with class balance)

ADHD- affected = 20094, Unaffected = 45803

METRIC	XGBOOST		XGBOOST BALANCED	
	MEAN	SD	MEAN	SD
Accuracy	0.710138	0.0051083	0.596397	0.00620875
Sensitivity	0.891732	0.00510639	0.585798	0.0104099
Specificity	0.296207	0.0131839	0.606998	0.0110402
F1 score	0.680386	0.00609839	0.596321	0.00620687
Precision	0.682649	0.00694207	0.596471	0.00621796
ROC-AUC Score	0.722084	0.00316626	0.635234	0.00931322

Filename: model\_2way\_adhd\_multidata.py & model\_2way\_adhd\_multidata\_balanced.py



# EXPERIMENT-19

**DATE: 13/03/24**

## Description and Input

Successfully implemented the 3-way classification model for classification of individuals into patients affected with ASD-only, ADHD-only or both using 4 datasets and XGBoost classifier - with and without Class balancing.

### DATASETS CHOSEN FOR MODEL:

1. df1 = basic\_medical\_screening-2023-07-21
2. df2 = background\_history\_child-2023-07-21
3. df3 = background\_history\_sibling-2023-07-21
4. df4 = individuals\_registration-2023-07-21


## Results

Number of X features = 28, Number of samples = 65897

ADHD- only= 2793, ASD- only = 29416, Both = 17301

Prediction	Value
0	ADHD only
1	ASD only
2	Both

METRIC	XGBOOST		XGBOOST BALANCED	
	MEAN	SD	MEAN	SD
Accuracy	0.637407	0.00591548	0.585153	0.0147246
Sensitivity	0.191192	0.023343	0.724324	0.0173023



Specificity	0.985872	0.00129488	0.797537	0.0198403
F1 score	0.624134	0.00557477	0.582287	0.0155735
Precision	0.62149	0.00572235	0.582548	0.0157496

Filename: model\_3way\_multidata.py & model\_3way\_multidata\_balanced.py

## EXPERIMENT-20

***DATE: 20/03/24***

### Description and Input

Obtained 2 sets of genomic data (.TSV file) on variant counts - ASD prediction in one and ADHD prediction in the other as the last column. This data contains numerical information about the number of each type of variant. I merged each of these datasets with the previous 4 datasets and counted the number of NaNs for features(columns) in this merged df - for each TSV file.

#### **DATASETS CHOSEN FOR MODEL:**

- 1) Server\_final\_updated\_data\_with\_target\_all\_variants + previous 4 datasets
- 2) Server\_final\_updated\_data\_with\_target\_all\_variants\_with\_adhd + previous 4 datasets

### Results

The results were stored in 2 excel sheets :

- 1) nan\_counts\_asd
- 2) nan\_counts\_adhd

Filename: dataset\_genomic\_features\_adhd.py and dataset\_genomic\_features\_asd.py

# EXPERIMENT-21

**DATE: 20/03/24**

## Description and Input

Successfully implemented the 2-way classification model for ASD and ADHD using 5 datasets and XGBoost classifier - with Class balancing.

### DATASETS CHOSEN FOR MODEL:

- 1) df1 = basic\_medical\_screening-2023-07-21
- 2) df2 = background\_history\_child-2023-07-21
- 3) df3 = background\_history\_sibling-2023-07-21
- 4) df4 = individuals\_registration-2023-07-21
- 5) df5 = server\_final\_updated\_data\_with\_target\_all\_variants.tsv

## Results

Number of X features = 43

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS - ASD

+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+		
Accuracy	0.842813	0.00702391
+-----+-----+-----+		
Sensitivity	0.848936	0.00850976
+-----+-----+-----+		
Specificity	0.836691	0.0117813
+-----+-----+-----+		
F1 score	0.842799	0.0070277
+-----+-----+-----+		
Precision	0.842942	0.00700073

```

+-----+-----+-----+
| ROC-AUC Score: | 0.842814 |      0.00702212 |
+-----+-----+-----+

```

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS - ADHD

```

+-----+-----+-----+
| Metric      | Mean | Standard Deviation |
+=====+=====+=====+
| Accuracy    | 0.649302 |      0.00659275 |
+-----+-----+-----+
| Sensitivity | 0.611879 |      0.0101125 |
+-----+-----+-----+
| Specificity | 0.686727 |      0.0158678 |
+-----+-----+-----+
| F1 score    | 0.648765 |      0.0064883 |
+-----+-----+-----+
| Precision   | 0.650235 |      0.00685656 |
+-----+-----+-----+
| ROC-AUC Score | 0.649303 |      0.00659524 |
+-----+-----+-----+

```

Filename: model\_2way\_asd\_genomic.py & model\_2way\_adhd\_genomic.py



## **EXTRA - NATURE PAPER FEATURES PRESENT IN BASIC MEDICAL SCREENING AND TRIED**

### Features

Preterm birth

Hypoxia at birth

Fetal alcohol syndrome - did not help alone

Bleeding into brain

Traumatic brain injury

Brain infection

Infection in pregnancy

Lead poisoning

## EXPERIMENT-22

**DATE: 28/03/24, 05/04/24**

### Description and Input

Successfully implemented the 2-way classification model using multiple datasets and XGBoost classifier - with Class balancing for classifying individuals into:

1. affected by ASD only v/s both ASD and ADHD
2. Affected by ADHD only v/s both ASD and ADHD

#### DATASETS CHOSEN FOR MODEL:

- 1) df1 = basic\_medical\_screening-2023-07-21
- 2) df2 = background\_history\_child-2023-07-21
- 3) df3 = background\_history\_sibling-2023-07-21
- 4) df4 = individuals\_registration-2023-07-21

### Results

#### CASE 1: ADHD ONLY V/S BOTH

The number of individuals affected is -> ADHD only: 2793, Both: 17301

#### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+		
Accuracy	0.725921	0.0180678
+-----+-----+-----+		
Sensitivity	0.740764	0.0294121
+-----+-----+-----+		
Specificity	0.711065	0.0225297
+-----+-----+-----+		
F1 score	0.725764	0.0180266

```

+-----+-----+-----+
| Precision   | 0.726458 | 0.0183136 |
+-----+-----+-----+
| ROC-AUC Score: | 0.797711 | 0.0185465 |
+-----+-----+-----+

```

Filename: model\_adhdonly\_both.py

### **CASE 2: ASD ONLY V/S BOTH**

The number of individuals affected is -> ASD only: 29416, Both: 17301

MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

```

+-----+-----+-----+
| Metric      | Mean | Standard Deviation |
+=====+=====+=====+
| Accuracy    | 0.660858 | 0.00533556 |
+-----+-----+-----+
| Sensitivity  | 0.59141 | 0.00780647 |
+-----+-----+-----+
| Specificity  | 0.730305 | 0.00795559 |
+-----+-----+-----+
| F1 score    | 0.659202 | 0.00535843 |
+-----+-----+-----+
| Precision   | 0.664047 | 0.00550446 |
+-----+-----+-----+
| ROC-AUC Score: | 0.715058 | 0.00476497 |
+-----+-----+-----+

```

Filename: model\_asdonly\_both.py



# EXPERIMENT-23

**DATE: 05/04/24**

## Description and Input

Successfully implemented the 2-way classification model removing dataset 3 and XGBoost classifier - with Class balancing for classifying individuals into:

1. Affected by ASD only v/s both ASD and ADHD
2. Affected by ADHD only v/s both ASD and ADHD -> This crashed as no of individuals with ADHD only became 0

### DATASETS CHOSEN FOR MODEL:

- 1) df1 = basic\_medical\_screening-2023-07-21
- 2) df2 = background\_history\_child-2023-07-21
- 3) df3 = background\_history\_sibling-2023-07-21 - REMOVED DATASET
- 4) df4 = individuals\_registration-2023-07-21

## Results

### ASD ONLY V/S BOTH


The number of individuals affected are:

ASD only: 29395

Both: 17301

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

Metric	Mean	Standard Deviation	
=====	=====	=====	=====
Accuracy	0.657419	0.00570264	
-----	-----	-----	-----
Sensitivity	0.578291	0.00892872	
-----	-----	-----	-----



Specificity	0.736547	0.00660572	
-------------	----------	------------	--

+-----+	+-----+	+-----+	+
---------	---------	---------	---

F1 score	0.655249	0.00582323	
----------	----------	------------	--

+-----+	+-----+	+-----+	+
---------	---------	---------	---

Precision	0.661479	0.00570259	
-----------	----------	------------	--

+-----+	+-----+	+-----+	+
---------	---------	---------	---

ROC-AUC Score:	0.707441	0.00707333	
----------------	----------	------------	--

+-----+	+-----+	+-----+	+
---------	---------	---------	---

Filename: model\_adhdonly\_both\_nodf3.py

# EXPERIMENT-24

**DATE: 05/04/24**

## Description and Input

Successfully implemented the 2-way classification model using the basic medical screening dataset and XGBoost classifier - with Class balancing for classifying individuals into:

1. Affected by ASD only v/s both ASD and ADHD
2. Affected by ADHD only v/s both ASD and ADHD

### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

## Results

### CASE1: ASD ONLY V/S BOTH

The number of individuals affected are:

ASD only: 49804

Both: 29827

Number of features = 11+5(nature paper) = 16, No of samples = 59654

### MEAN AND STANDARD DEVIATION OF VARIOUS METRICS

+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+		
Accuracy	0.546015	0.00541671
+-----+-----+-----+		
Sensitivity	0.672009	0.0152007
+-----+-----+-----+		
Specificity	0.420022	0.0132683
+-----+-----+-----+		

F1 score	0.538609	0.00540818
+-----+-----+-----+		
Precision	0.549202	0.00597462
+-----+-----+-----+		
ROC-AUC Score:	0.563649	0.00699528
+-----+-----+-----+		

Filename: model\_asdonly\_both\_basic\_medical.py

### **CASE2: ADHD ONLY V/S BOTH**

Sample size = 28852, Features = 11

The number of individuals affected are:

ADHD only: 14426, Both: 29827

### **MEAN AND STANDARD DEVIATION OF VARIOUS METRICS**

+-----+-----+-----+		
Metric	Mean	Standard Deviation
+=====+=====+=====+		
Accuracy	0.774332	0.00629648
+-----+-----+-----+		
Sensitivity	0.84944	0.00788158
+-----+-----+-----+		
Specificity	0.699226	0.0111601
+-----+-----+-----+		
F1 score	0.773037	0.00642595
+-----+-----+-----+		
Precision	0.780723	0.00613561
+-----+-----+-----+		
ROC-AUC Score:	0.837856	0.00623604
+-----+-----+-----+		

Filename: model\_adhdonly\_both\_basic\_medical.py

## EXPERIMENT-25

**DATE: 18/04/24**

### Description and Input

A repetition of exp 24 - classifying individuals into those Affected by ADHD only v/s both ASD and ADHD using 4 different classifiers and comparing results -> **THIS IS THE FINAL RESULT FOR POSTER**

#### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

### Results

Sample size = 28852, Features = 11

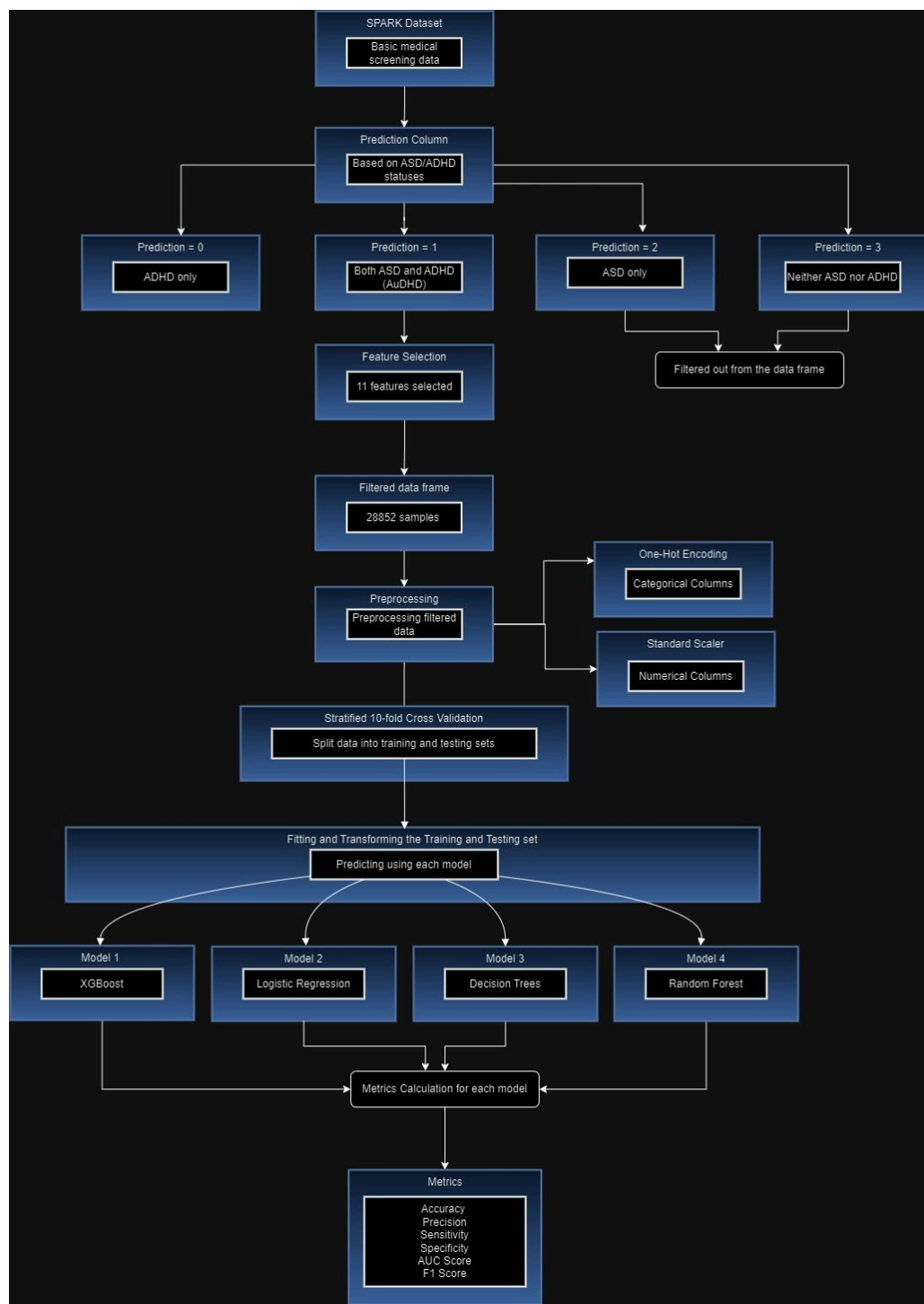
The number of individuals affected are:

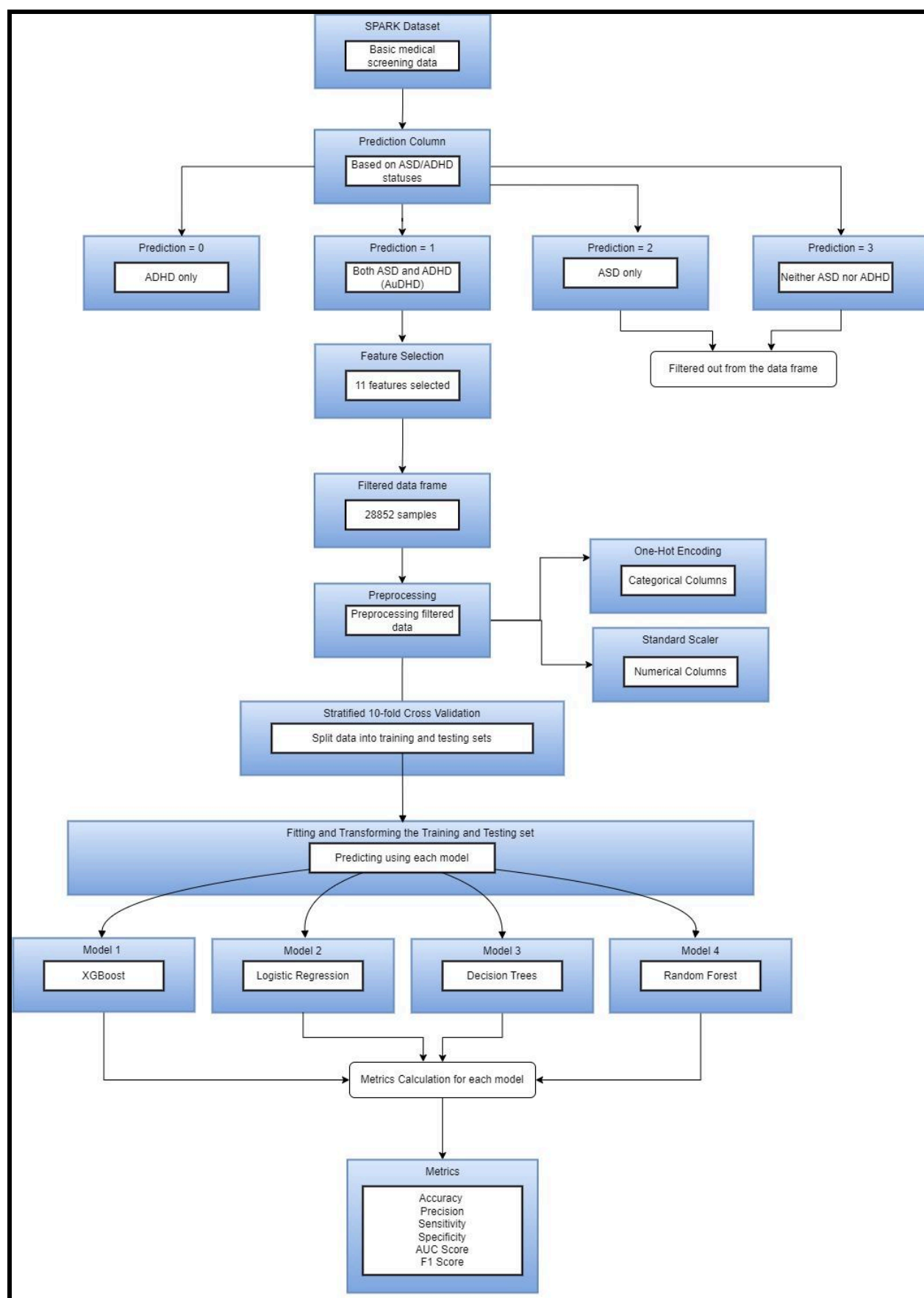
ADHD only: 14426, Both: 29827

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.774332	0.00629648	0.763449	0.00561958	0.769098	0.0066501	0.771871	0.00624375
Sensitivity	0.84944	0.00788158	0.872315	0.00733369	0.852351	0.00829936	0.850618	0.00836235
Specificity	0.699226	0.0111601	0.654584	0.0094651	0.685847	0.0110038	0.693125	0.0103141
F1 score	0.773037	0.00613561	0.760599	0.00577574	0.767472	0.00678412	0.770434	0.006342

Precision	0.7807 23	0.006 13561	0.776 608	0.005 7287	0.7768 27	0.006 6100 5	0.7788 43	0.006 26712
ROC-AUC Score	0.8378 56	0.006 23604	0.776 608	0.006 80111	0.8293 93	0.006 2120 5	0.8340 9	0.005 98536

## FLOWCHART FOR WORKFLOW







## EXPERIMENT-26

**DATE:** 01/05/24

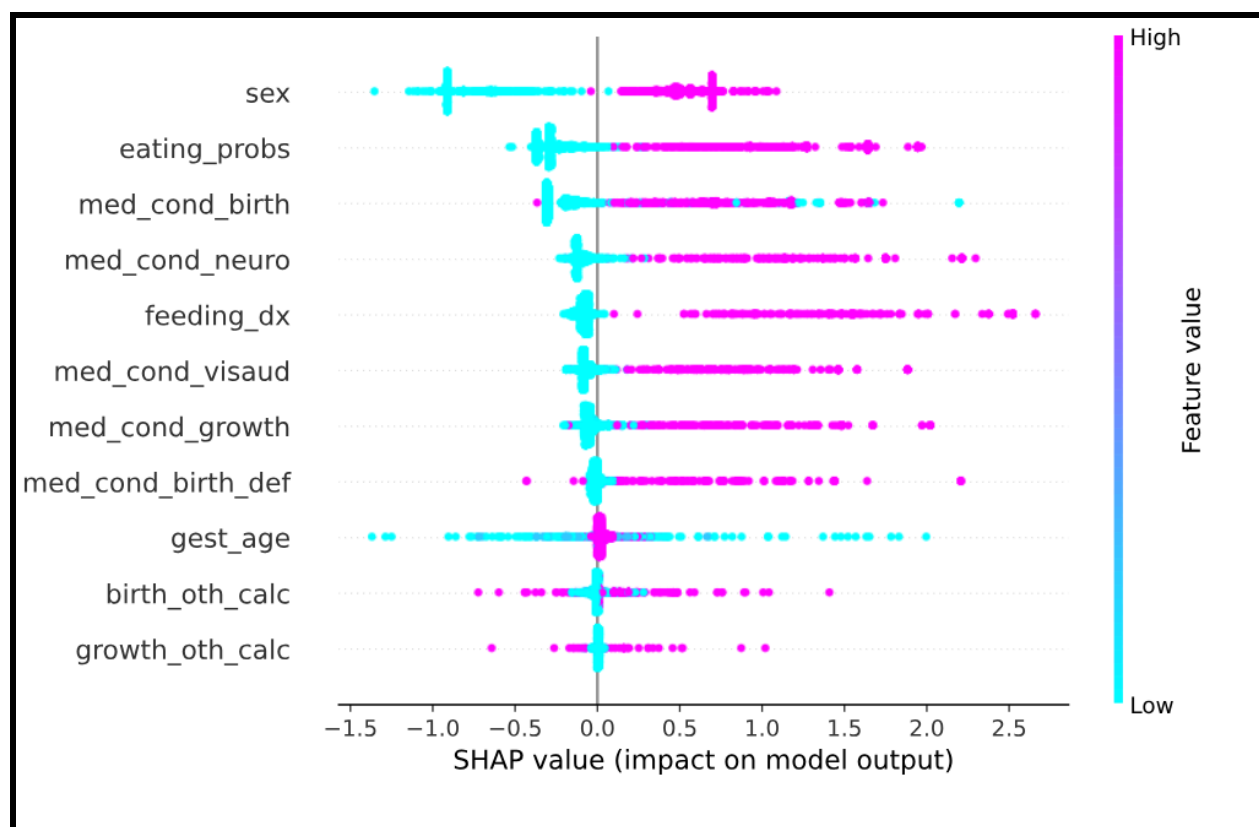
### Description and Input

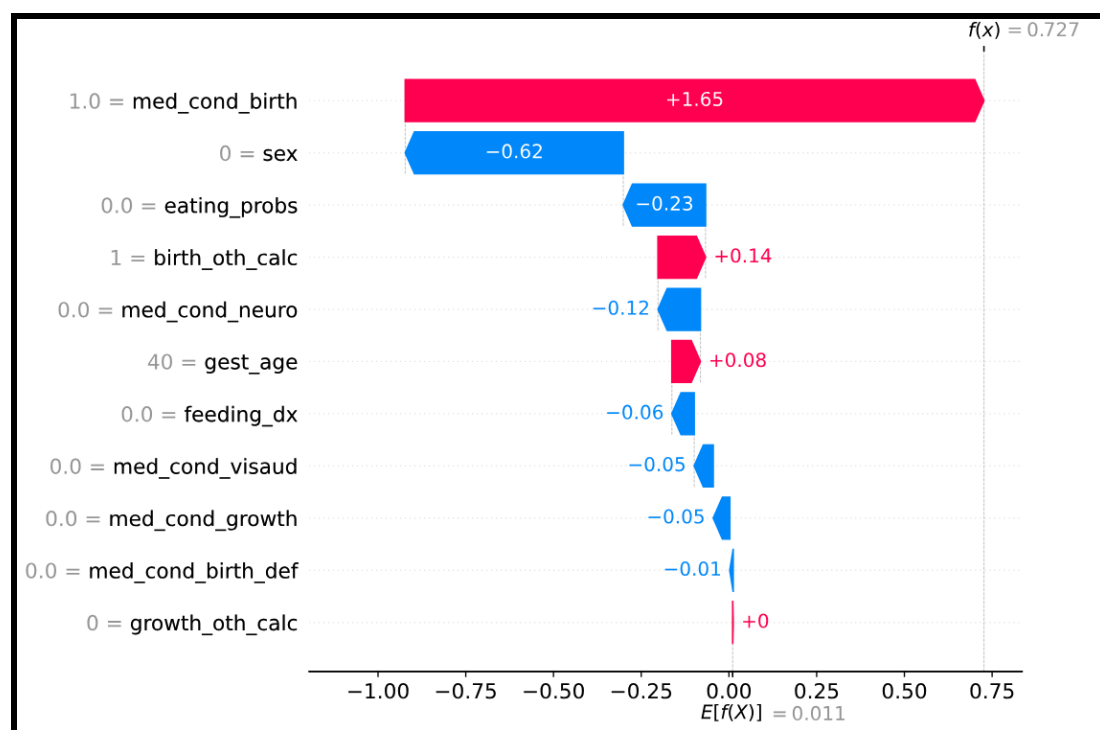
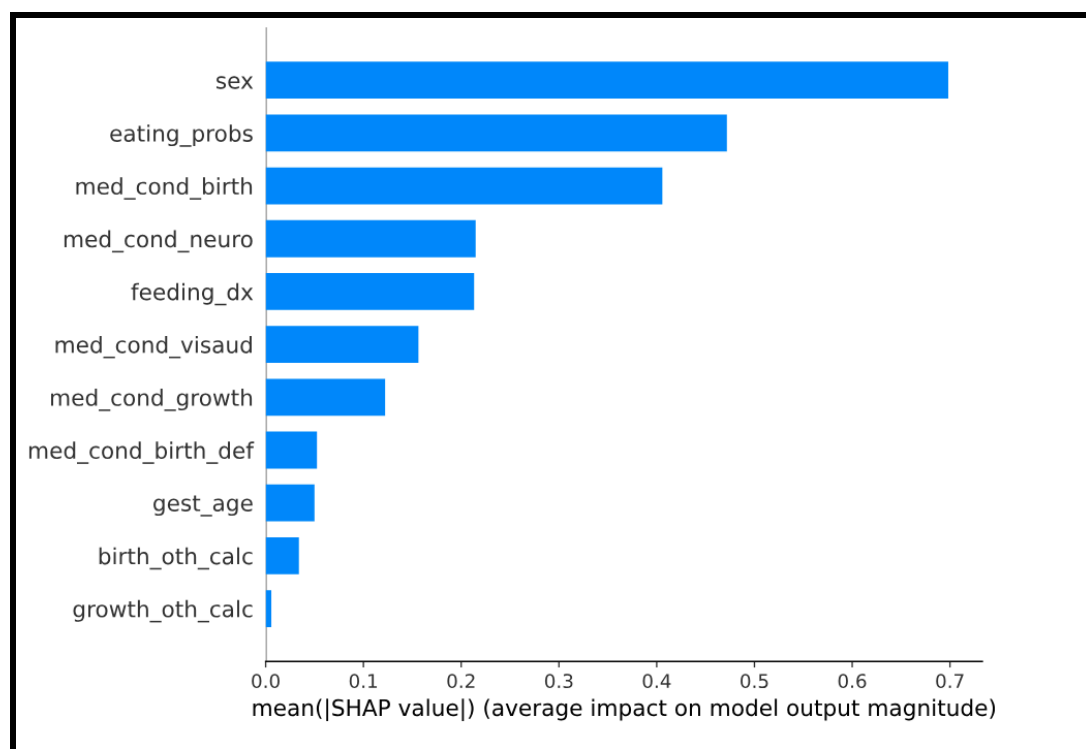
The model used in experiment 24 (XGBoost classifier) was subjected to SHAP analysis after preprocessing the encoded categorical columns. Different plots were obtained

#### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

### Results





Filename: model\_adhdonly\_both\_basic\_medical\_new.py

## EXPERIMENT-27

**DATE: 03/05/24**

### Description and Input

The model used in experiment 24 (ADHD only v/s AuDHD) was implemented without using sex as one of the features

#### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

### Results

Sample size = 28852, Features = 10

The number of individuals affected are:

ADHD only: 14426, Both: 29827

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.775233	0.00546441	0.76993	0.00648469	0.771732	0.00502226	0.773361	0.00563606
Sensitivity	0.847914	0.00814457	0.853044	0.00851792	0.850064	0.00809217	0.848954	0.00824876
Specificity	0.702553	0.00963431	0.686818	0.0144099	0.693402	0.00827227	0.69777	0.0097316
F1 score	0.774027	0.00553714	0.768301	0.00674435	0.770313	0.00506351	0.772046	0.00571314

Precision	0.7812 33	0.005 51508	0.777 702	0.005 9043	0.7786 22	0.005 2570 3	0.7798 15	0.005 69779
ROC-AUC Score	0.7895 1	0.006 6396	0.789 977	0.006 32248	0.7841 12	0.005 7639 5	0.7863 97	0.006 87734

Filename: model\_adhdonly\_both\_basic\_medical\_2.py

# EXPERIMENT-28

**DATE: 04/05/24**

## Description and Input

The plots for model used in experiment 24 (ADHD only v/s AuDHD) was implemented

### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

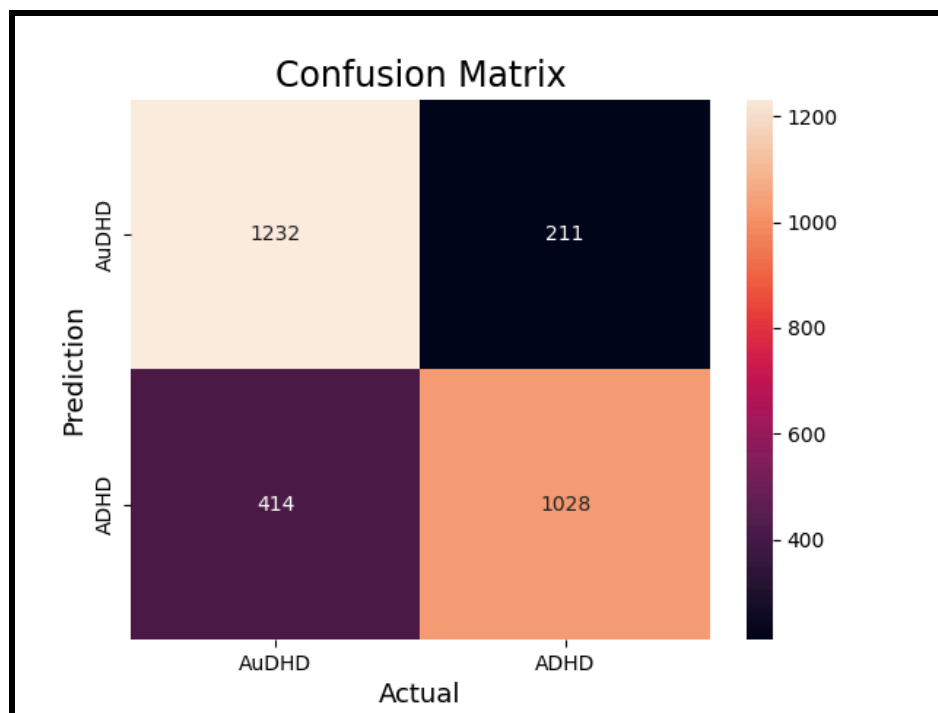
## Results

Sample size = 28852, Features = 11

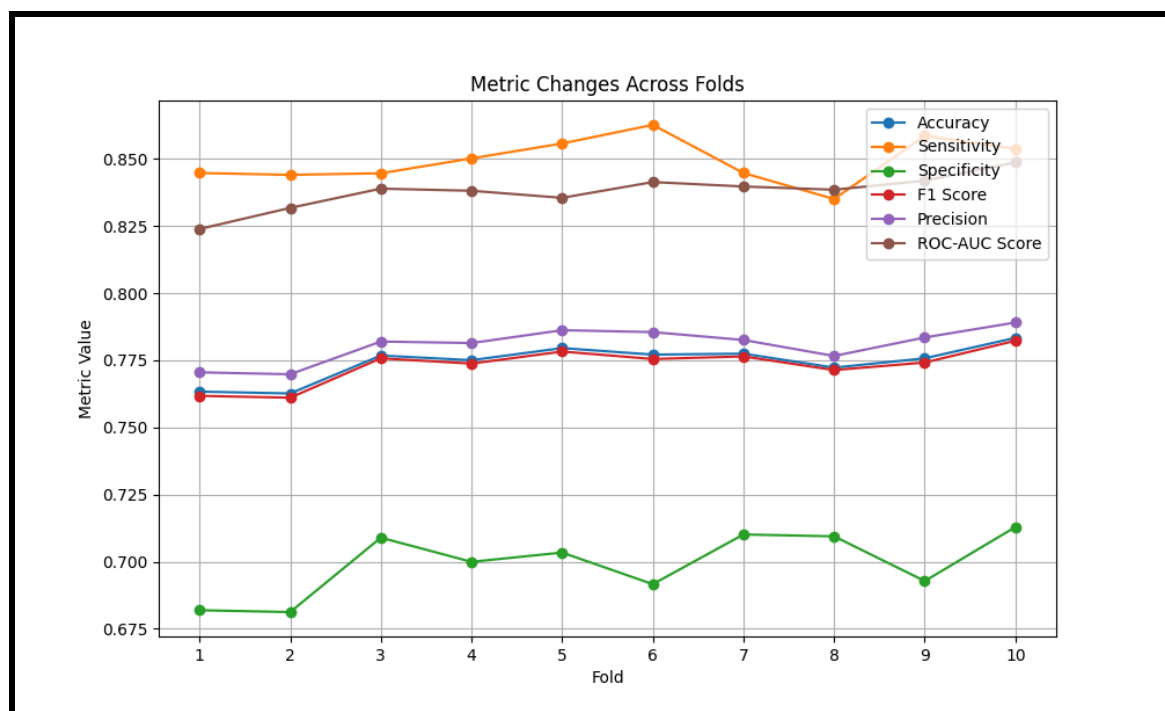
The number of individuals affected are:

ADHD only: 14426, Both: 29827

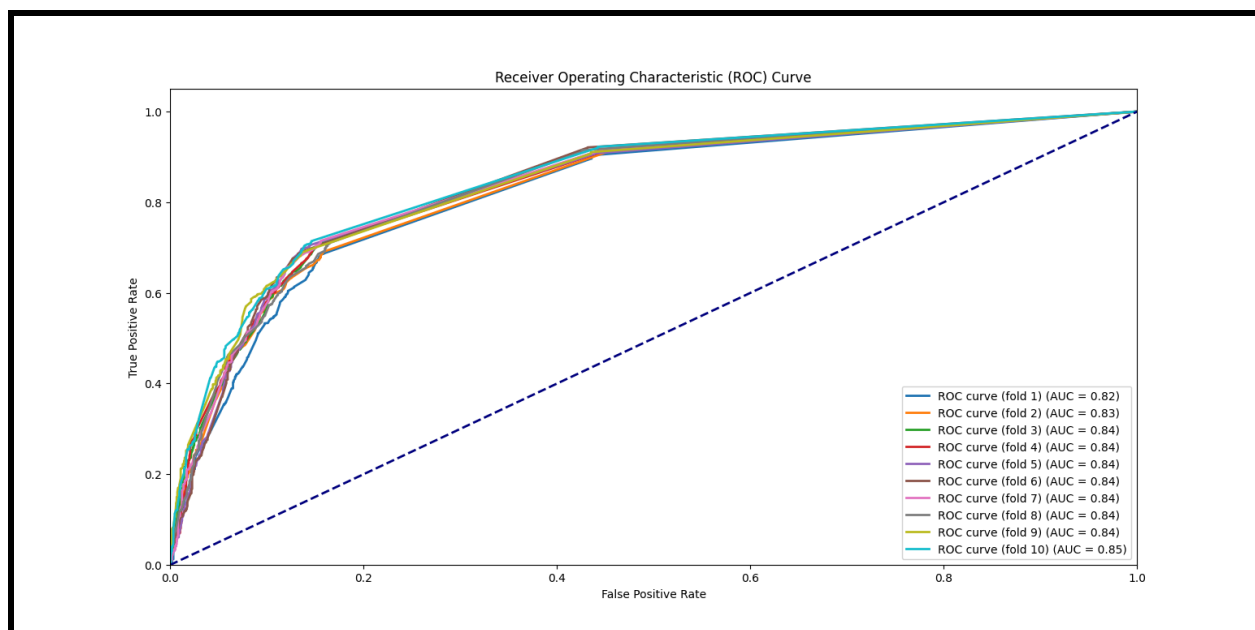
### Confusion matrix plot



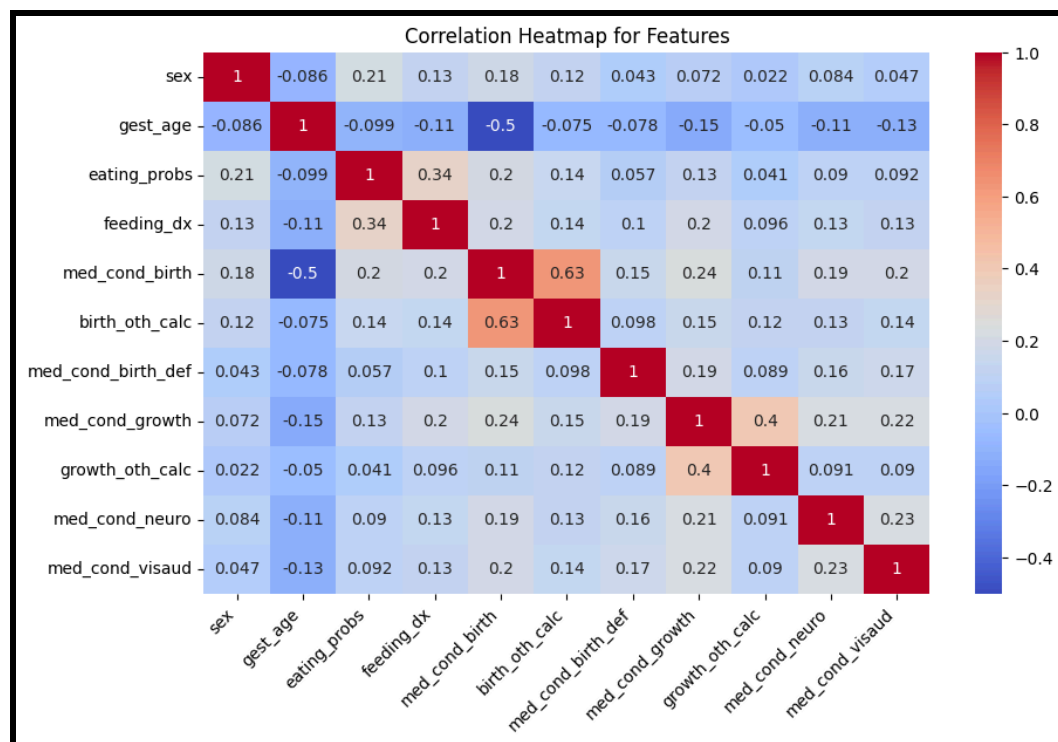
## Metrics fold changes



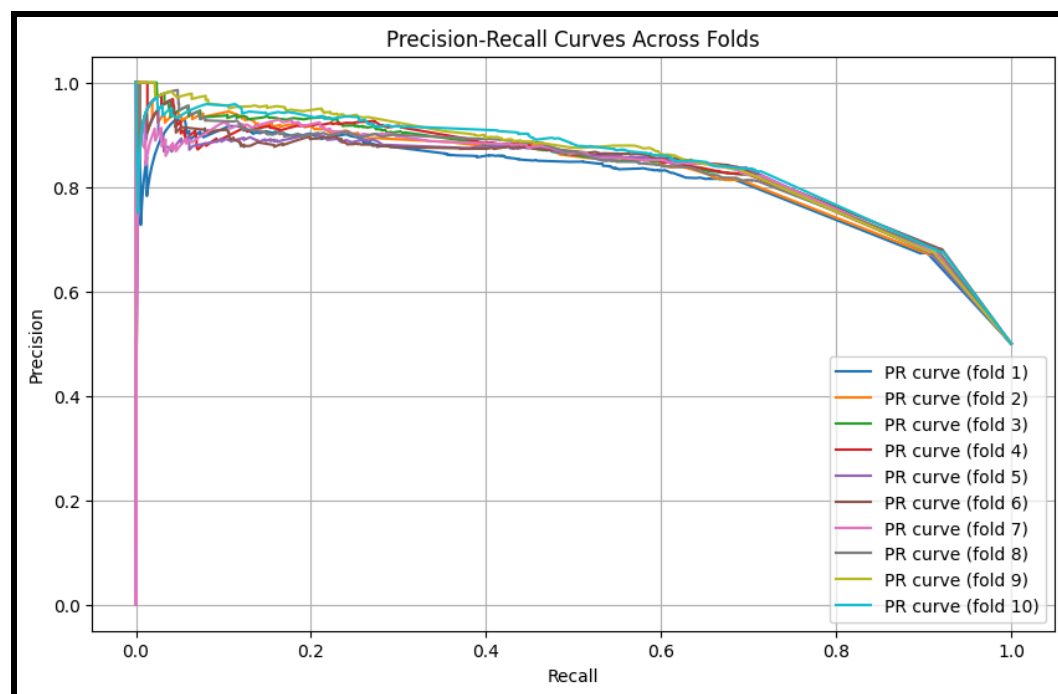
## ROC-Curve(All 10 folds)



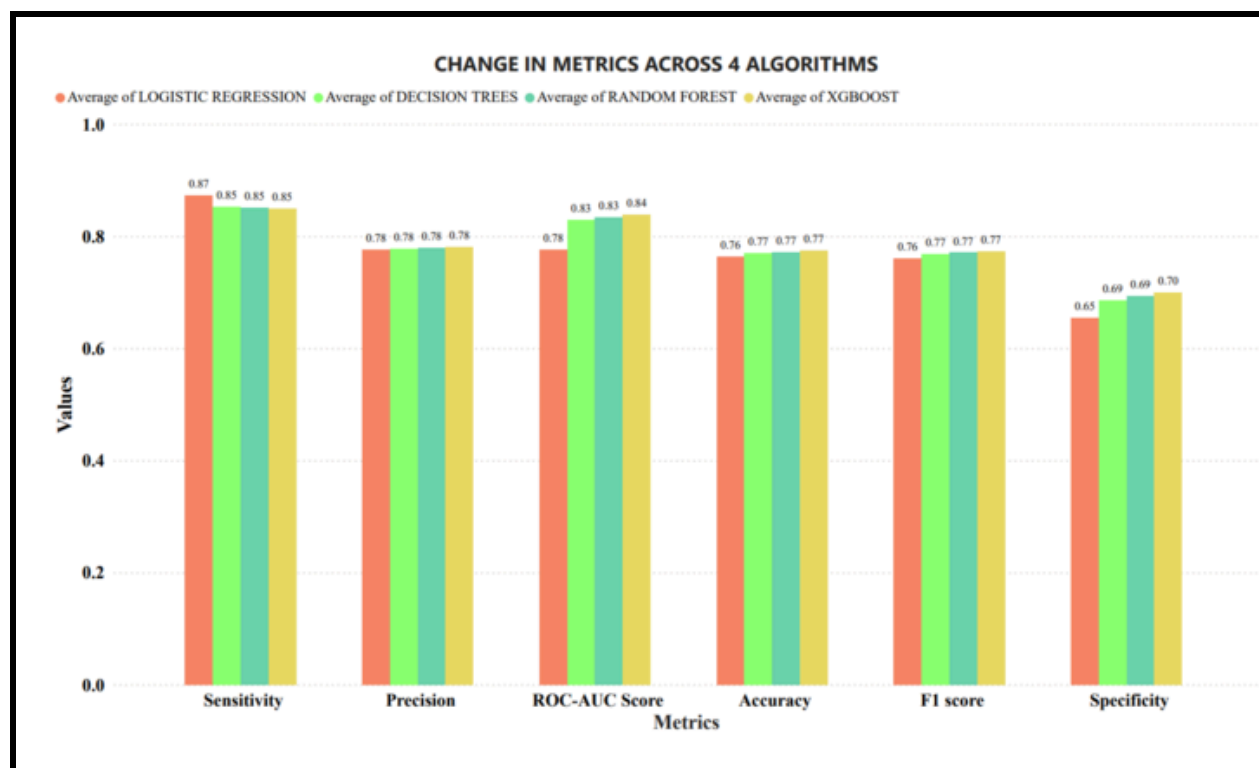
## Feature Correlation Heatmap



## Precision-Recall curve



### Metrics Change Across the 4 Classifiers



Filename: model\_adhdonly\_both\_basic\_medical\_new.py



## EXPERIMENT-29

**DATE: 05/05/24**

### Description and Input

The model used in experiment 24 (ADHD only v/s AuDHD) was implemented by class balancing based on sex instead of AuDHD/ADHD

#### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

### Results

The number of males in the dataset are: 27382

The number of females in the dataset are: 16871

The number of individuals affected are:

ADHD only: 14426, Both: 29827

Sample size = 33742, Features = 11

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.794736	0.00697491	0.794558	0.00638615	0.790617	0.00742432	0.792869	0.00725836
Sensitivity	0.652253	0.0112078	0.652577	0.0114871	0.654597	0.0117495	0.653951	0.0113246
Specificity	0.877246	0.00634492	0.876778	0.0054623	0.869384	0.0066536	0.873315	0.00647984
F1 score	0.791138	0.00711084	0.790991	0.00658026	0.787449	0.0075520	0.789507	0.00738632

						6		
Precision	0.7918 42	0.007 2634	0.791 651	0.006 65449	0.7875 96	0.007 6891 7	0.7899	0.007 53867
ROC-AUC Score	0.8398 79	0.006 26137	0.835 232	0.006 68988	0.8348 19	0.006 7662 1	0.8371 27	0.006 93749

Filename: model\_adhdonly\_both\_basic\_medical\_sex\_classbalance.py

## EXPERIMENT-30

**DATE: 06/05/24**

### Description and Input

The model used in experiment 24 (ADHD only v/s AuDHD) was implemented by age stratification. 3 sets of experiments were done (i.e. 3 age groups):

1. 0-2 years ( includes age 0 and 1 years)
2. 2-6 years ( includes age 2 to 5 years)
3. 6-10 years ( inclues age 6 to 9 years)

#### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

### Results

#### CASE 1: 0-2 YEARS ----> NOT RELIABLE

The number of individuals affected are:

ADHD only: 19

Both: 11

Sample size = 22, Features = 11

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.5166 67	0.283 333	0.733 333	0.351 188				
Sensitivity	0.65	0.45	0.8	0.4				
Specificity	0.35	0.45	0.7	0.458 258				

F1 score	0.4333 33	0.3	0.683 333	0.397 562				
Precision	0.3916 67	0.316 338	0.661 111	0.420 317				
ROC-AUC Score	0.6	0.435 89	0.7	0.458 258				

### **CASE 2: 2-6 YEARS**

The number of individuals affected are:

ADHD only: 254

Both: 2712

Sample size = 508, Features = 11

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.6278 04	0.059 2887	0.627 961	0.059 7866	0.6104 31	0.064 5274	0.5866 27	0.063 8246
Sensitivity	0.6926 15	0.047 5223	0.672 615	0.100 531	0.6850 77	0.060 5929	0.6490 77	0.118 875
Specificity	0.5633 85	0.111 669	0.582 769	0.099 2082	0.5356 92	0.115 115	0.5238 46	0.083 7275
F1 score	0.6246 27	0.061 317	0.624 582	0.061 0932	0.6063 01	0.066 41	0.5824 42	0.063 1818
Precision	0.6315 81	0.060 5925	0.631 822	0.059 2133	0.6142 26	0.064 8531	0.5914 23	0.067 7295
ROC-AUC Score	0.6209 69	0.066 0103	0.657 526	0.070 4169	0.5837 14	0.084 9531	0.5966 71	0.077 3416

### **CASE 3: 6-10 YEARS**

The number of individuals affected are:

ADHD only: 1295

Both: 8433

Sample size = 2590, Features = 11

METRIC	XGBOOST		LOGISTIC REGRESSION		DECISION TREES		RANDOM FOREST	
	MEAN	SD	MEAN	SD	MEAN	SD	MEAN	SD
Accuracy	0.643629	0.0351775	0.647104	0.0354372	0.620077	0.0245531	0.644402	0.0244466
Sensitivity	0.752153	0.0365042	0.742934	0.040305	0.729004	0.0278377	0.716577	0.0285023
Specificity	0.535075	0.052298	0.551366	0.0633908	0.511187	0.0310315	0.572194	0.053503
F1 score	0.638991	0.0365761	0.643172	0.0366321	0.615403	0.0251177	0.642005	0.0257016
Precision	0.650905	0.0354967	0.653502	0.0349924	0.626202	0.0258833	0.647987	0.0231872
ROC-AUC Score	0.683581	0.0333613	0.702239	0.0353258	0.660552	0.0208956	0.684112	0.0246874

Filename: model\_adhdonly\_both\_basic\_medical\_age\_strat.py

# EXPERIMENT-31

**DATE: 08/05/24**

## Description and Input

The model used in experiment 24 (ADHD only v/s AuDHD) was subjected to various types of statistical analysis

### DATASET CHOSEN FOR MODEL:

basic\_medical\_screening-2023-07-21

## Results

### **1. Confusion Matrix Analysis:**

Chi-Squared Test:

Chi-Squared Statistic: 942.9356814557027

P-value: 4.553633162985258e-207

Fisher's Exact Test:

Odds Ratio: 14.498431668841725

P-value: 1.0934394995396334e-221

*Interpretation:*

Based on the results of the statistical tests:

### 1. Chi-Squared Test:

- Chi-Squared Statistic: 942.94

- P-value: 4.55e-207

The chi-squared statistic is a measure of the discrepancy between the observed and expected frequencies in the contingency table. A larger chi-squared statistic indicates a greater discrepancy. The extremely small p-value (close to zero) suggests strong evidence against the null hypothesis, indicating that there is a significant association between the predicted and actual class labels.

## 2. Fisher's Exact Test:

- Odds Ratio: 14.50
- P-value: 1.09e-221

Fisher's exact test is another method for determining the association between categorical variables, particularly when the sample size is small. The odds ratio measures the strength of association between the two variables. In this case, the odds ratio of 14.50 indicates a strong association. The very small p-value (close to zero) suggests strong evidence against the null hypothesis, indicating that there is a significant association between the predicted and actual class labels.

In summary, both tests provide strong evidence to reject the null hypothesis, indicating a significant association between the predicted and actual class labels. This suggests that the model's predictions are not random and are indeed associated with the true class labels.