

Title: Autism Screening Adult Data Set

Informations :

Number of Instances:704

Attribute Characteristics: Integer

Number of Attributes:21

Date Donated 2017-12-24

Associated Tasks: Classification

Missing Values? Yes

Number of Web Hits: 84051

Dataset Content :

Attribute	Domain
1. A1_Score	{0,1}
2. A2_Score	{0,1}
3. A3_Score	{0,1}
4. A4_Score	{0,1}
5. A5_Score	{0,1}
6. A6_Score	{0,1}
7. A7_Score	{0,1}
8. A8_Score	{0,1}
9. A9_Score	{0,1}
10. A10_Score	{0,1}
11. age	numeric
12. gender	{f,m}
13. ethnicity	{White-European, Latino, Others, Black, Asian, 'Middle Eastern', Pasifika, 'South asian', Hispanic, Turkish, others}
14. jundice	{no, yes}
15. austin	{no, yes}
16. contry_of_res	{ 'United States', Brazil, Spain, Egypt, 'New Zealand', Bahamas, Burundi, Austria, Argentina, Jordan, Ireland, 'United Arab Emirates', Afghanistan, Lebanon, 'United Kingdom', 'South Africa', Italy, Pakistan, Bangladesh, Chile, France, China, Australia, Canada, Arabia', Netherlands, Romania, Sweden, Tonga, Oman, India, Philippines, 'Sri Lanka', 'Sierra Leone', Ethiopia, 'Viet Nam', Iran, 'Costa Rica', Germany, Mexico, Russia, Armenia, Iceland, Nicaragua, 'Hong Kong', Japan, Ukraine, Kazakhstan, American Samoa, Uruguay, Serbia, Portugal, M Republic', Cyprus}
17. used_app_before	{no, yes}
18. result	numeric
19. age_desc	{ '18 and more'}
20. relation	{Self, Parent, 'Health care

professional',Relative,Others}

21.Class/ASD

{NO,YES}

Dataset Description

Feature : Description

index : The participant's ID number

AX_Score: Score based on the Autism Spectrum Quotient (AQ) 10

item screening tool AQ-10

age : Age in years

gender : Male or Female

ethnicity: Ethnicities in text form

jaundice : Whether or not the participant was born with

jaundice?

austsm : Whether or not anyone in the immediate family has been diagnosed with autism?

country_of_res : Countries in text format

used_app_before : Whether the participant has used a screening app

result Score from the AQ-10 screening tool

age_desc : Age as categorical

relation : Relation of person who completed the test

Class/ASD : Participant classification

Importing Libraries

In [258...

```
# For dataframe and visualization
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from scipy.io import arff

# Processing data
from sklearn import preprocessing
from sklearn.preprocessing import StandardScaler

# Prepare Data for classification
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score

# Classification
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn import metrics
from sklearn.metrics import mean_absolute_error

# Comparing Classification
from sklearn.metrics import RocCurveDisplay
```

Reading Data

In [259...

```
# Load dataset
```

```
dataset = pd.read_table('Autism-Adult-Data.arff', sep = ',')
```

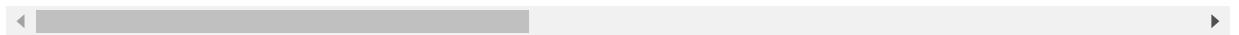
```
In [260... df = dataset.copy()
```

```
In [261... # Rename columns
df.columns = ['A1_Score', 'A2_Score', 'A3_Score', 'A4_Score', 'A5_Score', 'A6_Score', 'A7_Score', 'A8_Score', 'A9_Score', 'A10_Score', 'age', 'gender', 'ethnicity', 'jundice', 'austim', 'contry_of_res', 'used_app_before', 'resultnumeric', 'age_desc', 'relation']
```

```
In [262... df.head()
```

```
Out[262...
   A1_Score  A2_Score  A3_Score  A4_Score  A5_Score  A6_Score  A7_Score  A8_Score  A9_Score  A10_Score  age  gender  ethnicity  jundice  austim  contry_of_res  used_app_before  resultnumeric  age_desc  relation
0         1         1         0         1         0         0         0         1         1         1         1         1         1         0         1         1         0         1         1         1
1         1         1         0         1         1         0         1         1         1         1         1         1         1         0         1         1         0         1         1         1
2         1         1         0         1         0         0         1         1         1         1         1         1         1         0         1         1         0         1         1         1
3         1         0         0         0         0         0         0         0         1         1         1         1         1         0         1         1         0         1         1         1
4         1         1         1         1         1         0         1         1         1         1         1         1         1         0         1         1         0         1         1         1
```

5 rows × 21 columns



```
In [263... df.index #Describe index
```

```
Out[263... RangeIndex(start=0, stop=703, step=1)
```

```
In [264... df.shape
```

```
Out[264... (703, 21)
```

```
In [265... df.count() #Number of non-NA values
```

```
Out[265...
A1_Score      703
A2_Score      703
A3_Score      703
A4_Score      703
A5_Score      703
A6_Score      703
A7_Score      703
A8_Score      703
A9_Score      703
A10_Score     703
age           703
gender        703
ethnicity     703
jundice       703
austim        703
contry_of_res 703
used_app_before 703
resultnumeric 703
age_desc      703
relation      703
```

Class/ASD 703
dtype: int64

Feature Engineering

In [266...

```
df.info() #Info on DataFrame
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 703 entries, 0 to 702
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   A1_Score               703 non-null   int64
1   A2_Score               703 non-null   int64
2   A3_Score               703 non-null   int64
3   A4_Score               703 non-null   int64
4   A5_Score               703 non-null   int64
5   A6_Score               703 non-null   int64
6   A7_Score               703 non-null   int64
7   A8_Score               703 non-null   int64
8   A9_Score               703 non-null   int64
9   A10_Score              703 non-null   int64
10  age                   703 non-null   object
11  gender                 703 non-null   object
12  ethnicity              703 non-null   object
13  jundice                 703 non-null   object
14  austim                 703 non-null   object
15  contry_of_res          703 non-null   object
16  used_app_before        703 non-null   object
17  resultnumeric          703 non-null   int64
18  age_desc               703 non-null   object
19  relation               703 non-null   object
20  Class/ASD              703 non-null   object
dtypes: int64(11), object(10)
memory usage: 115.5+ KB
```

Some columns are object and some of them has string Yes or No, we need to replace them to boolean (0, 1)

In [267...

```
# Replace columns with number
df = df.replace("yes", 1)
df = df.replace("no", 0)
df = df.replace("YES", 1)
df = df.replace("NO", 0)
df = df.replace("f", 1)
df = df.replace("m", 0)
```

In [268...

```
df.info() #Info on DataFrame
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 703 entries, 0 to 702
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   A1_Score               703 non-null   int64
1   A2_Score               703 non-null   int64
2   A3_Score               703 non-null   int64
3   A4_Score               703 non-null   int64
4   A5_Score               703 non-null   int64
5   A6_Score               703 non-null   int64
6   A7_Score               703 non-null   int64
7   A8_Score               703 non-null   int64
8   A9_Score               703 non-null   int64
9   A10_Score              703 non-null   int64
```

```

10 age                703 non-null    object
11 gender              703 non-null    int64
12 ethnicity           703 non-null    object
13 jundice              703 non-null    int64
14 austim              703 non-null    int64
15 contry_of_res       703 non-null    object
16 used_app_before     703 non-null    int64
17 resultnumeric       703 non-null    int64
18 age_desc            703 non-null    object
19 relation            703 non-null    object
20 Class/ASD           703 non-null    int64
dtypes: int64(16), object(5)
memory usage: 115.5+ KB

```

In [269...

```

# Show missing values
MissingValues = {col:df[df[col] == "?"].shape[0] for col in df.columns}
MissingValues

```

Out[269...

```

{'A1_Score': 0,
 'A2_Score': 0,
 'A3_Score': 0,
 'A4_Score': 0,
 'A5_Score': 0,
 'A6_Score': 0,
 'A7_Score': 0,
 'A8_Score': 0,
 'A9_Score': 0,
 'A10_Score': 0,
 'age': 2,
 'gender': 0,
 'ethnicity': 95,
 'jundice': 0,
 'austim': 0,
 'contry_of_res': 0,
 'used_app_before': 0,
 'resultnumeric': 0,
 'age_desc': 0,
 'relation': 95,
 'Class/ASD': 0}

```

Replace '?' values of Age by mean

In [270...

```

# Replace '?' values by NaN
for j in range(df.shape[0]):
    if(df.iloc[j,10]=='?'):
        df.iloc[j,10]=np.NaN

```

In [271...

```

# Drop NaN value
df.dropna(inplace=True)

```

In [272...

```

# replace dtype object to int
df["age"] = df["age"].astype(str).astype(int)

```

In [273...

```
df["age"].describe()
```

Out[273...

```

count    701.000000
mean      29.703281
std       16.518660
min       17.000000
25%       21.000000
50%       27.000000

```

```
75%      35.000000
max      383.000000
Name: age, dtype: float64
```

```
In [274... # dropping record number 51
df.iloc[51,10]=np.NaN
```

```
In [275... df[df['age'] == df['age'].max()]['age']
```

```
Out[275... 6      64.0
Name: age, dtype: float64
```

```
In [276... df["age"].describe()
```

```
Out[276... count      700.000000
mean       29.198571
std        9.717718
min        17.000000
25%        21.000000
50%        27.000000
75%        35.000000
max        64.000000
Name: age, dtype: float64
```

```
In [277... df['age'].fillna((df['age'].mean()), inplace=True)
```

Replace '?' values of ethnicity by 'Others' and 'others' by 'Others'

```
In [279... # There is values that are the same : '?', 'Others' and 'others'
df['ethnicity'].unique()
```

```
Out[279... array(['Latino', 'White-European', '?', 'Others', 'Black', 'Asian',
      "Middle Eastern '", 'Pasifika', "'South Asian'", 'Hispanic',
      'Turkish', 'others'], dtype=object)
```

```
In [280... # Replace '?' with 'others'
df['ethnicity'] = df['ethnicity'].replace('?', 'others')
```

```
In [281... # Replace '?' with 'Others'
df['ethnicity'] = df['ethnicity'].replace('others', 'Others')
```

```
In [282... # Every missing values are now as 'Others'
df['ethnicity'].unique()
```

```
Out[282... array(['Latino', 'White-European', 'Others', 'Black', 'Asian',
      "Middle Eastern '", 'Pasifika', "'South Asian'", 'Hispanic',
      'Turkish'], dtype=object)
```

Replace '?' values of relation by a mode of relation

```
In [283... # Here we only have '?' as missing values
df['relation'].unique()
```

```
Out[283... array(['Self', 'Parent', '?', "'Health care professional'", 'Relative',
```

```
'Others'], dtype=object)
```

```
In [284... # Replace the missing value with modal value of the columns
df['relation'] = df['relation'].replace('?', df['relation'].mode()[0])
```

```
In [285... # Show results
df['relation'].unique()
```

```
Out[285... array(['Self', 'Parent', "'Health care professional'", 'Relative',
      'Others'], dtype=object)
```

```
In [286... # No more missing values !
df.isnull().sum() #Number of NA values
```

```
Out[286... A1_Score      0
A2_Score      0
A3_Score      0
A4_Score      0
A5_Score      0
A6_Score      0
A7_Score      0
A8_Score      0
A9_Score      0
A10_Score     0
age           0
gender        0
ethnicity     0
jundice       0
austim        0
contry_of_res 0
used_app_before 0
resultnumeric 0
age_desc     0
relation      0
Class/ASD     0
dtype: int64
```

```
In [287... # Now every columns has the right type
df.info() #Info on DataFrame
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 701 entries, 0 to 702
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   A1_Score              701 non-null   int64
1   A2_Score              701 non-null   int64
2   A3_Score              701 non-null   int64
3   A4_Score              701 non-null   int64
4   A5_Score              701 non-null   int64
5   A6_Score              701 non-null   int64
6   A7_Score              701 non-null   int64
7   A8_Score              701 non-null   int64
8   A9_Score              701 non-null   int64
9   A10_Score             701 non-null   int64
10  age                   701 non-null   float64
11  gender                701 non-null   int64
12  ethnicity             701 non-null   object
13  jundice               701 non-null   int64
14  austim                701 non-null   int64
15  contry_of_res         701 non-null   object
16  used_app_before       701 non-null   int64
17  resultnumeric         701 non-null   int64
18  age_desc              701 non-null   object
```

```

19 relation          701 non-null    object
20 Class/ASD         701 non-null    int64
dtypes: float64(1), int64(16), object(4)
memory usage: 120.5+ KB

```

: number of functions in the data framework
 Column: Features header in the Dataframe
 Non-null Count: Counter of nonzero values for each Dataframe
 function
 Type: type of data stored for each function of the data frame

Summary

In [288...

```
df.describe() #Statistical summary of DataFrame
```

Out[288...

	A1_Score	A2_Score	A3_Score	A4_Score	A5_Score	A6_Score	A7_Score	A8_
count	701.000000	701.000000	701.000000	701.000000	701.000000	701.000000	701.000000	701.0
mean	0.723252	0.452211	0.457917	0.496434	0.499287	0.285307	0.416548	0.6
std	0.447710	0.498066	0.498582	0.500344	0.500357	0.451883	0.493339	0.4
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
50%	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.0
75%	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.0
max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.0

count: number of examples counted for the selected function
 mean: arithmetic mean for the selected function
 std: standard deviation for the selected function
 min: minimum value presented by the examples for the selected function
 25%: first quartile calculated on the examples for the selected function
 50%: second quartile calculated on the examples for the selected function
 75%: third quartile calculated on examples for selected feature
 max: maximum value presented by the examples for the selected function

In [289...

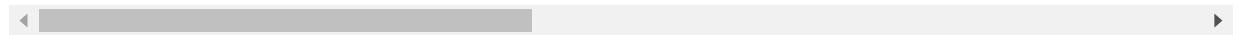
```
df.head()
```

Out[289...

	A1_Score	A2_Score	A3_Score	A4_Score	A5_Score	A6_Score	A7_Score	A8_Score	A9_Scc
0	1	1	0	1	0	0	0	1	
1	1	1	0	1	1	0	1	1	
2	1	1	0	1	0	0	1	1	
3	1	0	0	0	0	0	0	1	

A1_Score	A2_Score	A3_Score	A4_Score	A5_Score	A6_Score	A7_Score	A8_Score	A9_Score
4	1	1	1	1	1	0	1	1

5 rows × 21 columns



Visualization

In [313...

```
# Let's see the diversity of autism
print(df['gender'].value_counts())
men = df.value_counts(["gender"])[0]
women = df.value_counts(["gender"])[1]

name = ['men', 'women']
data = [men, women]
plt.title("Patients with autism by gender", fontsize = 15)

plt.pie(data, labels=name, startangle=90, shadow=True)
plt.axis('equal')
plt.show()
```

0 366

1 335

Name: gender, dtype: int64

Patients with autism by gender

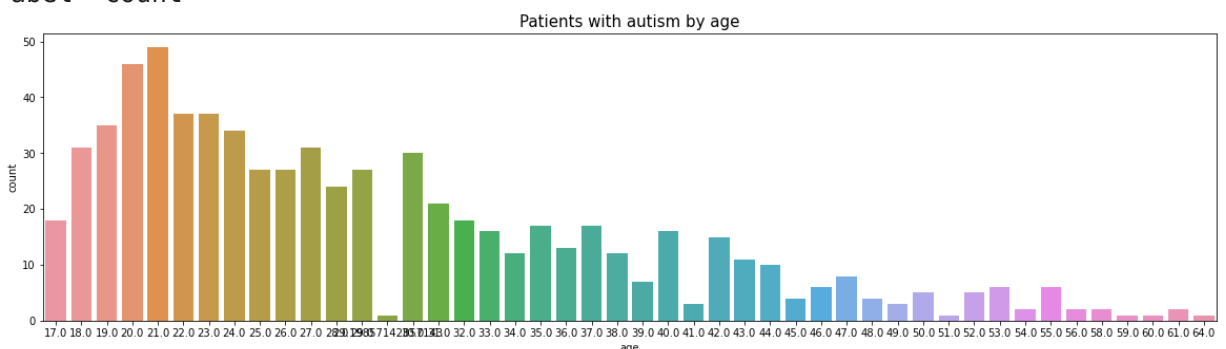


In [314...

```
plt.figure(figsize=(20,5))
plt.title("Patients with autism by age", fontsize = 15)
sns.countplot(x="age",data=df)
```

Out[314...

<AxesSubplot:title={'center': 'Patients with autism by age'}, xlabel='age', ylabel='count'>



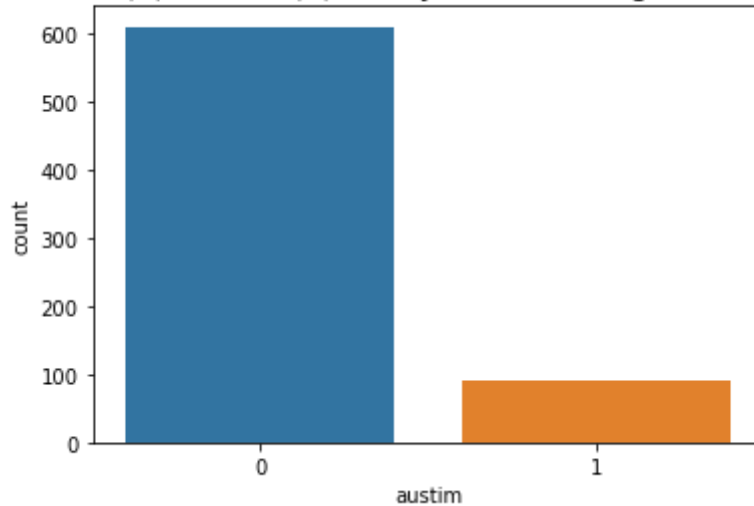
```
In [293... # Let's see if Patients has family member diagnosed with autism
print(df['austim'].value_counts())
sns.countplot(x="austim", data=df)
plt.title("Patients with(1) and not(0) family member diagnosed with autism",
plt.show()
```

```
0    610
```

```
1     91
```

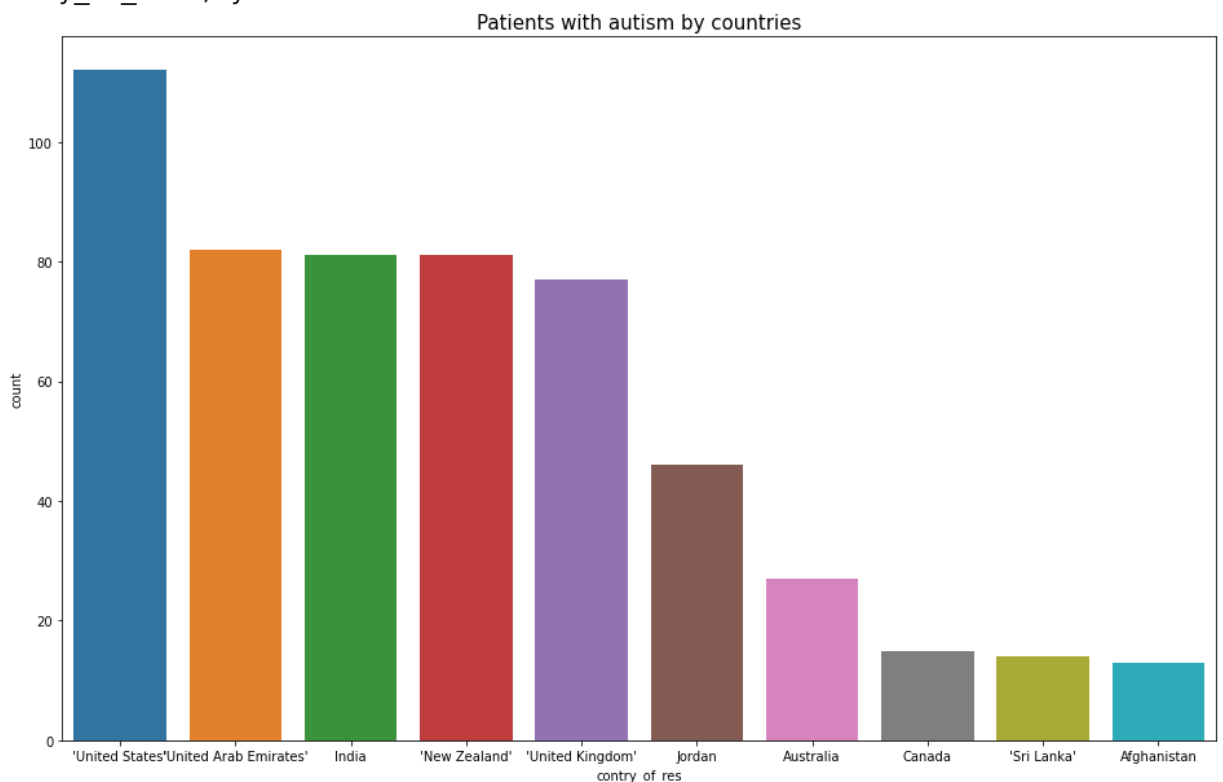
```
Name: austim, dtype: int64
```

Patients with(1) and not(0) family member diagnosed with autism



```
In [294... # The top 10 countries with autism
plt.figure(figsize = (16, 10))
plt.title("Patients with autism by countries", fontsize = 15)
plt.xlabel('countries')
plt.ylabel('number of austim')
order=df["contry_of_res"].value_counts().nlargest(10).index
# plt.bar(df.value_counts(["contry_of_res"])[1],df['austim'])
sns.countplot(x="contry_of_res", data=df, order=order)
```

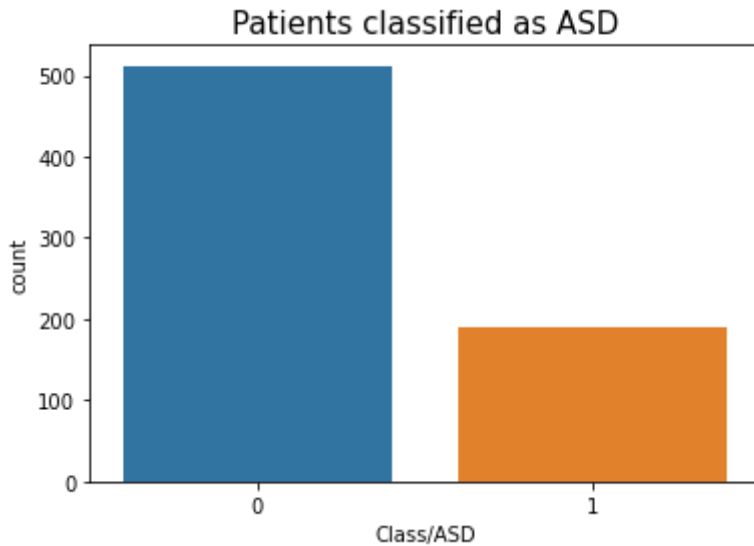
```
Out[294... <AxesSubplot:title={'center': 'Patients with autism by countries'}, xlabel='co
ntry_of_res', ylabel='count'>
```



In [291...

```
plt.title("Patients classified as ASD", fontsize = 15)
sns.countplot(x="Class/ASD", data=df)
print(df["Class/ASD"].describe())
```

```
count    701.000000
mean      0.269615
std       0.444077
min       0.000000
25%      0.000000
50%      0.000000
75%      1.000000
max       1.000000
Name: Class/ASD, dtype: float64
```

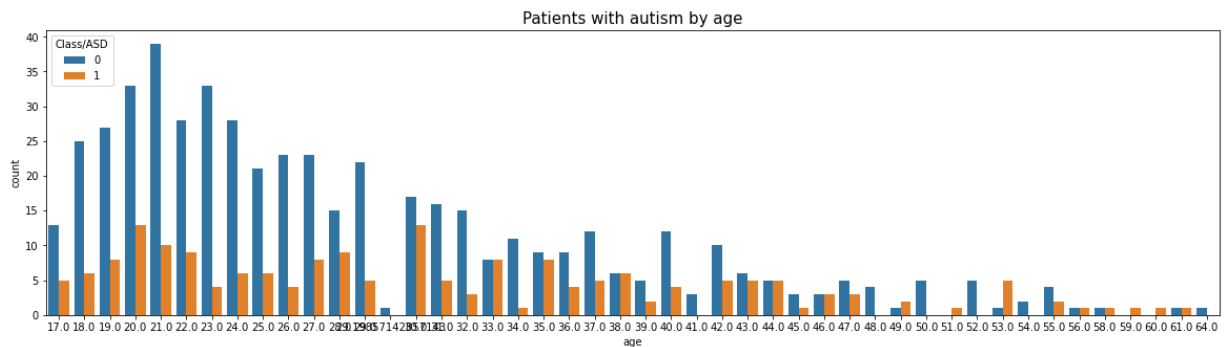


In [317...

```
plt.figure(figsize=(20,5))
plt.title("Patients with autism by age", fontsize = 15)
sns.countplot(x="age", hue="Class/ASD", data=df)
```

Out[317...

```
<AxesSubplot:title={'center': 'Patients with autism by age'}, xlabel='age', ylabel='count'>
```

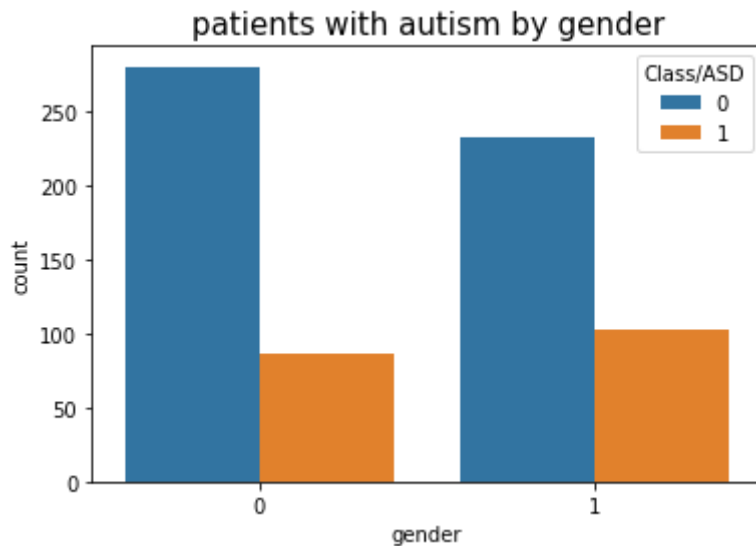


In [316...

```
plt.title("Patients with autism by gender", fontsize = 15)
sns.countplot(x="gender", hue="Class/ASD", data=df)
```

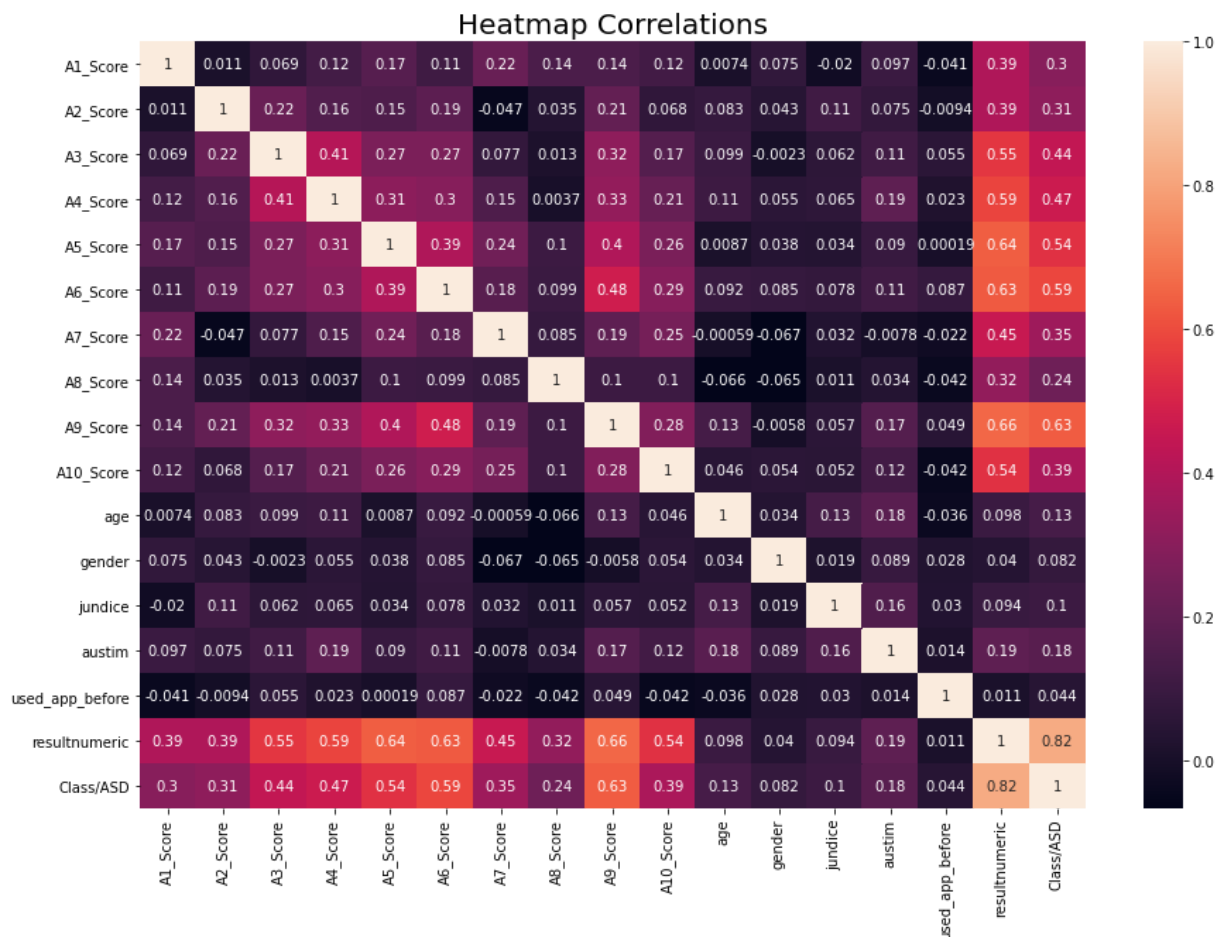
Out[316...

```
<AxesSubplot:title={'center': 'patients with autism by gender'}, xlabel='gender', ylabel='count'>
```



In [295...

```
# Correlation between dataset columns
plt.figure(figsize = (15, 10))
plt.title("Heatmap Correlations", fontsize = 20)
sns.heatmap(df.corr(), annot = True)
plt.show()
```



Pré-processing

In [296...

```
# Dropp Unwanted columns
df.drop(['age_desc'], axis = 1, inplace = True)
```

Split the data

```
In [297... X = df.drop("Class/ASD", axis = 1) # select all other feature except "Class"
y = df['Class/ASD']
```

Due to the presence of data expressed with different location, normalization must be performed by using the `get_dummies()` method.

```
In [298... X = pd.get_dummies(X)
```

The data need to be split in training set and testing set

```
In [299... X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.8)
```

```
In [300... print(f"X = {X.shape}")
print(f"Y = {y.shape}")
```

```
X = (701, 98)
Y = (701,)
```

```
In [301... print(f"X_train = {X_train.shape}")
print(f"Y_train = {y_train.shape}\n")
print(f"X_test = {X_test.shape}")
print(f"Y_test = {y_test.shape}")
```

```
X_train = (140, 98)
Y_train = (140,)
```

```
X_test = (561, 98)
Y_test = (561,)
```

Support Vector Classification

```
In [302... # Apply SVC
svc = SVC(random_state=3)
svc.fit(X_train, y_train)
pred_svc = svc.predict(X_test)
print(classification_report(y_test, pred_svc))
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, pred_svc))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, pred_svc))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test,
```

	precision	recall	f1-score	support
0	0.92	1.00	0.96	414
1	1.00	0.74	0.85	147
accuracy			0.93	561
macro avg	0.96	0.87	0.90	561
weighted avg	0.94	0.93	0.93	561

```
Mean Absolute Error: 0.0677361853832442
Mean Squared Error: 0.0677361853832442
Root Mean Squared Error: 0.2602617631986001
```

```
In [303... A = np.array([ round(metrics.precision_score(y_test, pred_svc),4),
round(metrics.recall_score(y_test, pred_svc),4),
round(metrics.f1_score(y_test, pred_svc),4)])
A = np.reshape(A, (1, 3))
A
```

Out[303... array([[1. , 0.7415, 0.8516]])

Random Forest Classifier

```
In [304... # Apply RFC
rfc = RandomForestClassifier(random_state=3)
rfc.fit(X_train, y_train)
pred_RFR = rfc.predict(X_test)
print(classification_report(y_test, pred_RFR))
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, pred_RFR))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, pred_RFR))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test,
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	414
1	1.00	0.99	0.99	147
accuracy			1.00	561
macro avg	1.00	0.99	1.00	561
weighted avg	1.00	1.00	1.00	561

Mean Absolute Error: 0.0035650623885918
Mean Squared Error: 0.0035650623885918
Root Mean Squared Error: 0.05970814340265321

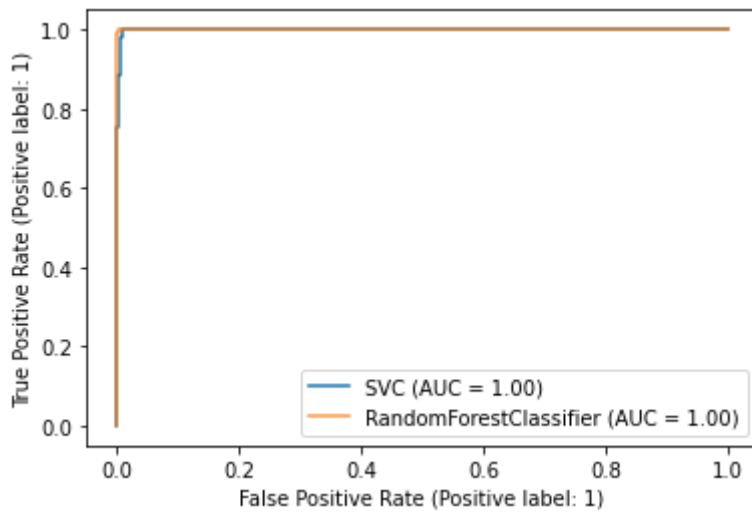
```
In [305... B = np.array([ round(metrics.precision_score(y_test, pred_RFR),4),
               round(metrics.recall_score(y_test, pred_RFR),4),
               round(metrics.f1_score(y_test, pred_RFR),4)])
B = np.reshape(B, (1, 3))
B
```

Out[305... array([[1. , 0.9864, 0.9932]])

Classification Comparison

Roc Curve

```
In [306... svc_disp = RocCurveDisplay.from_estimator(svc, X_test, y_test)
ax = plt.gca()
rfc_disp = RocCurveDisplay.from_estimator(rfc, X_test, y_test, ax=ax, alpha=0.5)
regressor_disp = ()
plt.show()
```

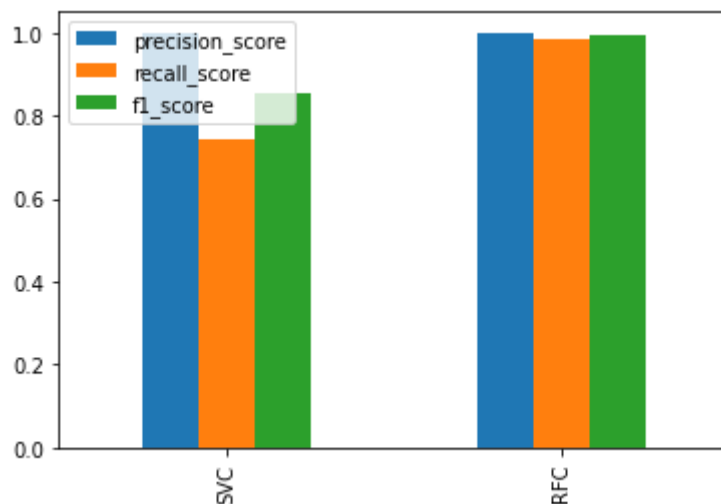


We cannot clearly see the difference with the Roc curve even if we can briefly see that RFC is slightly more eccentric than an SVC.

Plot Bar

In [307...

```
# plot scoring to see the difference
Data = np.reshape([A, B], (2, 3))
fig = pd.DataFrame(Data, columns=["precision_score", "recall_score", "f1_score"])
fig.plot.bar();
plt.show()
print([A, B])
```



```
[array([[1.      , 0.7415, 0.8516]]), array([[1.      , 0.9864, 0.9932]])]
```

We will use Random Forest Classifier. He has the best scores and is the most eccentric curve in the ROC curve.