# CIND-820 Capstone Project: An ML Tool to Detect Heart Disease

- Robert M. Pineau
- 941-049-371
- Supervisor: Dr. Ceni Babaoglu

**Install Required Modules:**

```
1 #!pip install matplotlib
2 #!pip install graphviz
```

**Load Required Libraries:**

```
 1 import sys
 2 from google.colab import drive
 3 import math
 4 from statistics import mean, stdev
 5 import pandas as pd
 6 import numpy as np
 7 from scipy import stats
 8 import plotly
 9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 from sklearn import tree
13 from sklearn.naive_bayes import MultinomialNB
14 from sklearn.naive_bayes import GaussianNB
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.svm import SVC
17 from sklearn.tree import DecisionTreeClassifier
18 from sklearn.neighbors import KNeighborsClassifier
19 from sklearn.ensemble import RandomForestClassifier
20 from xgboost import XGBClassifier
21
22 from sklearn.metrics import confusion_matrix
23 from sklearn.model_selection import train_test_split
24 from sklearn.metrics import classification_report
25 from sklearn import metrics
26 from sklearn.metrics import accuracy_score
27 from sklearn.metrics import f1_score
28
29 from sklearn.model_selection import KFold
30 from sklearn.model_selection import cross_validate
```

```
31 from sklearn.model_selection import cross_val_score
32 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
33 import imblearn
34 from imblearn.over_sampling import SMOTE
35
36 import graphviz
```

## Obtain Data-Set from Google Drive:

```
1 # Mounting google colab, this will prompt first time each session.
2 drive.mount('/content/drive',force_remount=True)
3 dataset_file = "/content/drive/My Drive/Colab Notebooks/heart_statlog_cleveland_hungary_fi
4 df=pd.read_csv(dataset_file,sep=',')
5 df.name = "Original Data-Set"
6 print(df.name)
7 df.head(3)
```

```
Mounted at /content/drive
Original Data-Set
```

| | age | sex | chest pain type | resting bp s | cholesterol | fasting blood sugar | resting ecg | max heart rate | exercise angina | oldpeak | sl |
|---|-----|-----|------|------|------|------|------|------|------|------|----|
| 0 | 40 | 1 | 2 | 140 | 289 | 0 | 0 | 172 | 0 | 0.0 | |
| 1 | 49 | 0 | 3 | 160 | 180 | 0 | 0 | 156 | 0 | 1.0 | |

## Clean up Column Names:

```
1 #Rename the columns to be nicer, no spaces.
2 df=df.rename(columns={"age": "Age", "sex": "Sex", "chest pain type": "ChestPainType", "res
3 df=df.rename(columns={"cholesterol":"Cholesterol","fasting blood sugar": "FastingBloodSuga
4 df=df.rename(columns={"chest pain type": "ChestPainType", "resting bp s": "RestingBP_s", "
5 df=df.rename(columns={"resting ecg": "RestingECG", "max heart rate": "MaxHeartRate", "exer
6 df=df.rename(columns={"oldpeak":"OldPeak", "ST slope": "ST_Slope", "target": "Target"})
7 df.name = "Original Data-Set"
8 df.head(3)
```

| | Age | Sex | ChestPainType | RestingBP_s | Cholesterol | FastingBloodSugar | RestingECG | Ma |
|---|-----|-----|------|------|------|------|------|----|
| 0 | 40 | 1 | 2 | 140 | 289 | 0 | 0 | |
| 1 | 49 | 0 | 3 | 160 | 180 | 0 | 0 | |
| 2 | 37 | 1 | 2 | 130 | 283 | 0 | 1 | |

## Datatypes and Quantities:

```
1 #Check data types.
2 print(df.dtypes)
```

```
Age                   int64
Sex                   int64
ChestPainType         int64
RestingBP_s           int64
Cholesterol           int64
FastingBloodSugar     int64
RestingECG            int64
MaxHeartRate          int64
ExerciseAngina        int64
OldPeak             float64
ST_Slope              int64
Target                int64
dtype: object
```

```
1 #Datatypes, counts, etc.
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1190 entries, 0 to 1189
Data columns (total 12 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   Age                1190 non-null   int64
 1   Sex                1190 non-null   int64
 2   ChestPainType      1190 non-null   int64
 3   RestingBP_s        1190 non-null   int64
 4   Cholesterol        1190 non-null   int64
 5   FastingBloodSugar  1190 non-null   int64
 6   RestingECG         1190 non-null   int64
 7   MaxHeartRate       1190 non-null   int64
 8   ExerciseAngina     1190 non-null   int64
 9   OldPeak            1190 non-null   float64
 10  ST_Slope           1190 non-null   int64
 11  Target             1190 non-null   int64
dtypes: float64(1), int64(11)
memory usage: 111.7 KB
```

## Check for Missing and NULL entries:

```
1 #Check for NULL or missing entries. (none)
2 print(df.isna().any())
3 print("\n\n")
4 print(df.isnull().any())
```

```
Age                 False
Sex                 False
ChestPainType       False
RestingBP_s         False
```

```
Cholesterol         False
FastingBloodSugar   False
RestingECG          False
MaxHeartRate        False
ExerciseAngina      False
OldPeak             False
ST_Slope            False
Target              False
dtype: bool
```

```
Age                 False
Sex                 False
ChestPainType       False
RestingBP_s         False
Cholesterol         False
FastingBloodSugar   False
RestingECG          False
MaxHeartRate        False
ExerciseAngina      False
OldPeak             False
ST_Slope            False
Target              False
dtype: bool
```

### Check for Duplicate Entries:

```python
1  #Look for rows that are 100% identical to each other.
2  #
3  dup_count = sum(df.duplicated())
4  print(f"There are {dup_count} duplicate rows in this dataset.\n")
5
6  #Droping any duplicate entries.
7  df = df.drop_duplicates(ignore_index = True)
8  df.name = "Original Data-Set"
9
10 df.info()
```

```
There are 272 duplicate rows in this dataset.

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   Age                918 non-null    int64
 1   Sex                918 non-null    int64
 2   ChestPainType      918 non-null    int64
 3   RestingBP_s        918 non-null    int64
 4   Cholesterol        918 non-null    int64
 5   FastingBloodSugar  918 non-null    int64
 6   RestingECG         918 non-null    int64
```

```
 7   MaxHeartRate        918 non-null    int64
 8   ExerciseAngina      918 non-null    int64
 9   OldPeak             918 non-null    float64
10   ST_Slope            918 non-null    int64
11   Target              918 non-null    int64
dtypes: float64(1), int64(11)
memory usage: 86.2 KB
```

## Check for Out-Of-Bound Entries for Nominal and Binary attributes:

```
1  #Check for out of bound entries(outliers) for nominal and binary attributes (including the
2  #Since all nominal and binary attributes have a valid contiguous integer range, ie 0-1, or
3  #we only need to look for those outside the range.
4  valid_values = {'Sex': [0,1], 'ChestPainType': [1,2,3,4], 'FastingBloodSugar': [0,1], 'Res
5
6  for col in ('Sex', 'ChestPainType', 'FastingBloodSugar', 'RestingECG', 'ExerciseAngina', '
7    valid = np.array(valid_values[col])
8    max_valid = valid.max()
9    min_valid = valid.min()
10   print(f"For attribute '{col}': Valid MAX: {max_valid}, Valid MIN: {min_valid}")
11   these_outliers = df[((df[col] < min_valid) | (df[col] > max_valid))]
12
13   if (these_outliers.shape[0] > 1):
14     print(f"For attribute '{col}': There are {these_outliers.shape[0]} outliers:\n")
15     print(these_outliers)
16     print("\n\n")
17   elif (these_outliers.shape[0] == 1):
18     print(f"For attribute '{col}': There is {these_outliers.shape[0]} outlier:\n")
19     print(these_outliers)
20     print("\n\n")
21   else:
22     print(f"For attribute '{col}': There are no outliers.\n")
```

```
For attribute 'Sex': Valid MAX: 1, Valid MIN: 0
For attribute 'Sex': There are no outliers.

For attribute 'ChestPainType': Valid MAX: 4, Valid MIN: 1
For attribute 'ChestPainType': There are no outliers.

For attribute 'FastingBloodSugar': Valid MAX: 1, Valid MIN: 0
For attribute 'FastingBloodSugar': There are no outliers.

For attribute 'RestingECG': Valid MAX: 2, Valid MIN: 0
For attribute 'RestingECG': There are no outliers.

For attribute 'ExerciseAngina': Valid MAX: 1, Valid MIN: 0
For attribute 'ExerciseAngina': There are no outliers.

For attribute 'ST_Slope': Valid MAX: 3, Valid MIN: 1
For attribute 'ST_Slope': There is 1 outlier:

     Age  Sex  ChestPainType  RestingBP_s  Cholesterol  FastingBloodSugar  \
```

```
516   68    1              3         150          195                      1

      RestingECG  MaxHeartRate  ExerciseAngina  OldPeak  ST_Slope  Target
516            0           132               0      0.0         0       1
```

```
For attribute 'Target': Valid MAX: 1, Valid MIN: 0
For attribute 'Target': There are no outliers.
```

### Remove Out-Of-Bound Entry:

```
 1 #From above, there is a problem with one entry regarding the ST_Slope attribute, it is zer
 2 #
 3 #The documentation at https://ieee-dataport.org/open-access/heart-disease-dataset-comprehe
 4 #Shows a range of 0-2, but in the definition of the mapped nominal values it shows:
 5 #
 6 # -- Value 1: upsloping
 7 # -- Value 2: flat
 8 # -- Value 3: downsloping
 9 #
10 print(df['ST_Slope'].value_counts().sort_index())
11 print("\n")
12 #
13 #
14 #Since there is only one entry out of range, making the assumption that the correct range
15 #
16 #Will simply drop this entry.
17 df = df[df['ST_Slope'] != 0]
18 df.name = "Original Data-Set"
19
20 df.info()
```

```
0      1
1    395
2    459
3     63
Name: ST_Slope, dtype: int64


<class 'pandas.core.frame.DataFrame'>
Int64Index: 917 entries, 0 to 917
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   Age              917 non-null    int64
 1   Sex              917 non-null    int64
 2   ChestPainType    917 non-null    int64
 3   RestingBP_s      917 non-null    int64
 4   Cholesterol      917 non-null    int64
 5   FastingBloodSugar  917 non-null  int64
```

```
 6   RestingECG           917 non-null    int64
 7   MaxHeartRate         917 non-null    int64
 8   ExerciseAngina       917 non-null    int64
 9   OldPeak              917 non-null    float64
10   ST_Slope             917 non-null    int64
11   Target               917 non-null    int64
dtypes: float64(1), int64(11)
memory usage: 93.1 KB
```

## Basic Statistics for All Attributes:

```
1 #Basic Statistics of the dataset.(Measures of Center/Central Tendency, and Measures of Var
2 df.describe().T
```

|  | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| Age | 917.0 | 53.495093 | 9.425601 | 28.0 | 47.0 | 54.0 | 60.0 | 77.0 |
| Sex | 917.0 | 0.789531 | 0.407864 | 0.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| ChestPainType | 917.0 | 3.251908 | 0.931502 | 1.0 | 3.0 | 4.0 | 4.0 | 4.0 |
| RestingBP_s | 917.0 | 132.377317 | 18.515114 | 0.0 | 120.0 | 130.0 | 140.0 | 200.0 |
| Cholesterol | 917.0 | 198.803708 | 109.443764 | 0.0 | 173.0 | 223.0 | 267.0 | 603.0 |
| FastingBloodSugar | 917.0 | 0.232279 | 0.422517 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| RestingECG | 917.0 | 0.604144 | 0.806161 | 0.0 | 0.0 | 0.0 | 1.0 | 2.0 |
| MaxHeartRate | 917.0 | 136.814613 | 25.473732 | 60.0 | 120.0 | 138.0 | 156.0 | 202.0 |
| ExerciseAngina | 917.0 | 0.404580 | 0.491078 | 0.0 | 0.0 | 0.0 | 1.0 | 1.0 |
| OldPeak | 917.0 | 0.888332 | 1.066749 | -2.6 | 0.0 | 0.6 | 1.5 | 6.2 |
| ST_Slope | 917.0 | 1.637950 | 0.607270 | 1.0 | 1.0 | 2.0 | 2.0 | 3.0 |
| Target | 917.0 | 0.552890 | 0.497466 | 0.0 | 0.0 | 1.0 | 1.0 | 1.0 |

## Visualizing All Attributes:

```
1 #Assigning descriptive labels for all possible values for all nominal/binary attributes.
2 #
3 labels = {'Sex': ['Female', 'Male'], 'ChestPainType': ['Typical Angina', 'A-Typical Angina
4          'FastingBloodSugar': ['<= 120 mg/dl', '> 120 mg/dl'], 'RestingECG': ['Normal', '
5          'ExerciseAngina':['No', 'Yes'], 'ST_Slope': ['Upsloping', 'Flat', 'Downsloping']
6
7
```
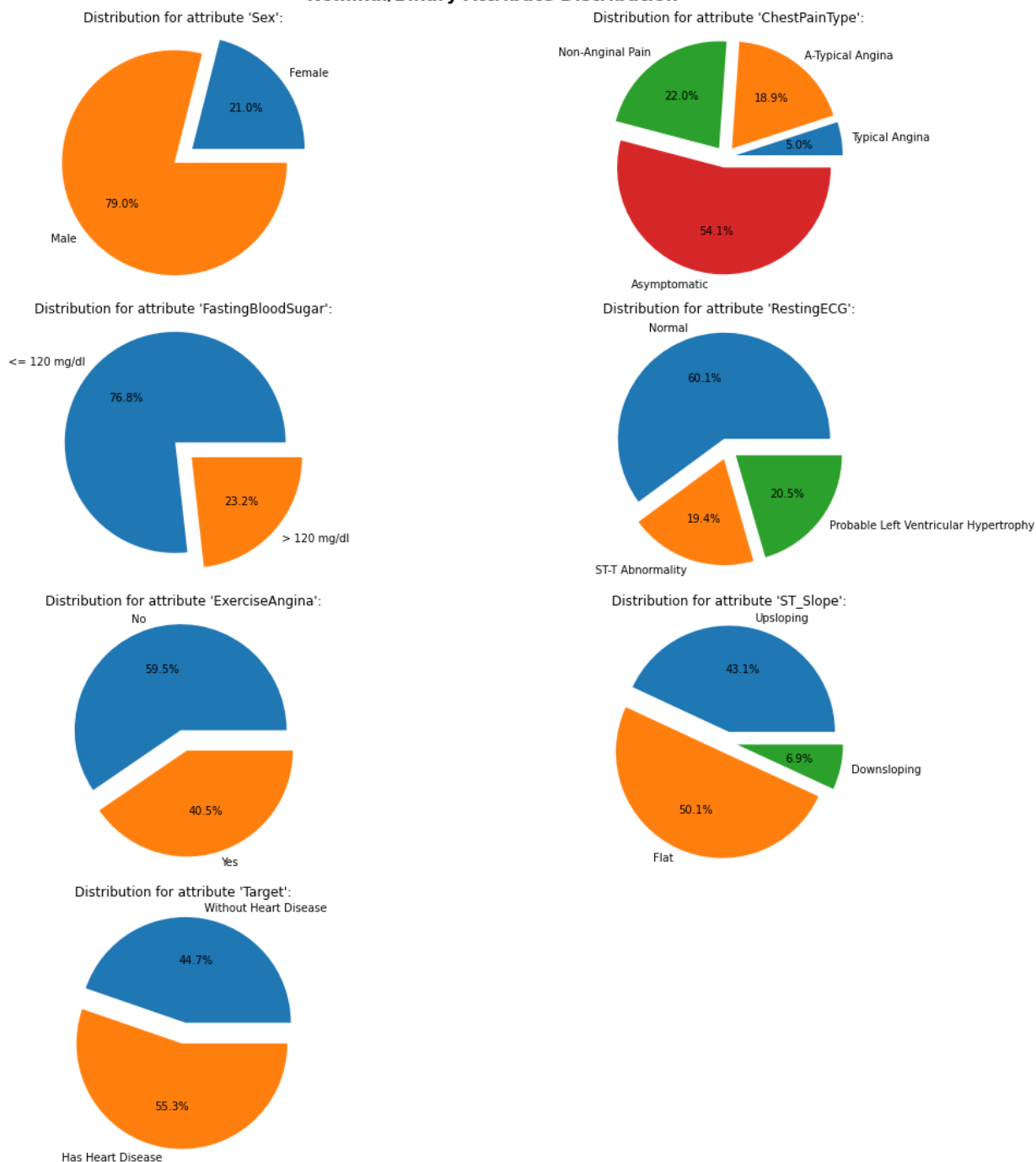
```
1 #Creating subsets of the data for a series of interesting plots to help with visualization
```

```python
2  #
3  #Breaking up the dataset into two groups, those with heart disease and those without.
4  with_heart_disease = df[df['Target'] == 1]
5  no_heart_disease = df[df['Target'] == 0]
6
7  #Breaking up the dataset into four groups, by Sex, and those with heart disease and those
8  with_heart_disease_male = df[(df['Target'] == 1) & (df['Sex'] == 1)]
9  with_heart_disease_female = df[(df['Target'] == 1) & (df['Sex'] == 0)]
10 no_heart_disease_male = df[(df['Target'] == 0) & (df['Sex'] == 1)]
11 no_heart_disease_female = df[(df['Target'] == 0) & (df['Sex'] == 0)]
12
13 #For these groups remove the "Sex" column from the data.
14 with_heart_disease_male = with_heart_disease_male.drop("Sex",axis=1)
15 with_heart_disease_female = with_heart_disease_female.drop("Sex",axis=1)
16 no_heart_disease_male = no_heart_disease_male.drop("Sex",axis=1)
17 no_heart_disease_female = no_heart_disease_female.drop("Sex",axis=1)
```

```python
1  #Visualizing the distribution of the categorical attributes, including the class variable.
2  ax=1
3  plt.figure(figsize=(15,15))
4
5  for col in ('Sex', 'ChestPainType', 'FastingBloodSugar', 'RestingECG', 'ExerciseAngina','S
6    plt.subplot(4,2,ax)
7    these_labels = labels[col]
8    plt.title(f"Distribution for attribute '{col}':")
9    plt.pie(df[col].value_counts().sort_index(),
10        autopct = '%1.1f%%', labels=these_labels,
11        explode=tuple([0.1] * len(these_labels)))
12   plt.axis('equal')
13   ax+=1
14
15 plt.suptitle('Nominal/Binary Attribute Distribution',y=1.01, size = 16, color = 'black', w
16 plt.tight_layout()
17 plt.savefig("nominal_dist.pdf",dpi=1200, bbox_inches='tight')
```
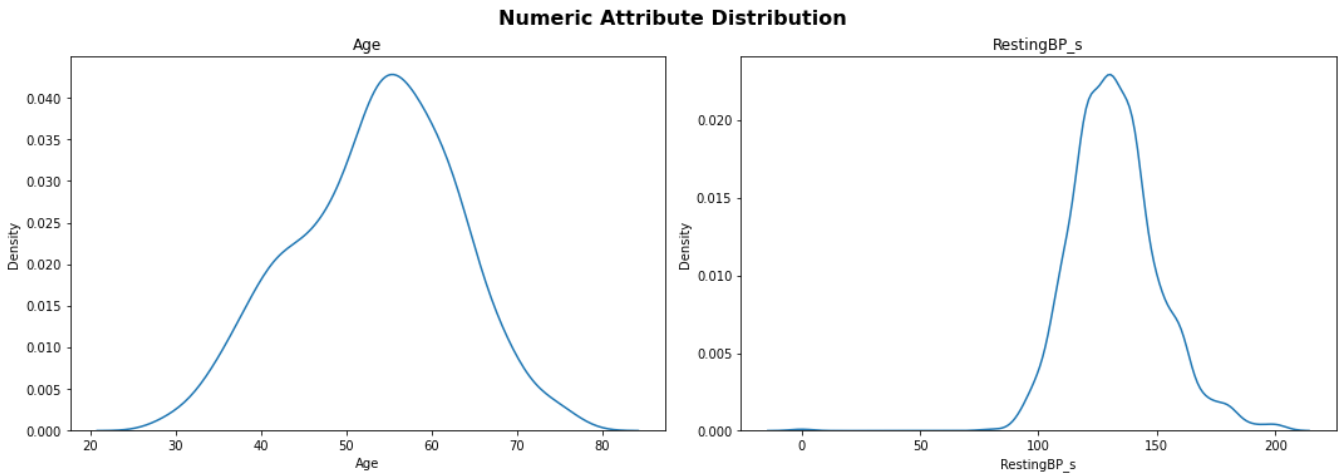
## Nominal/Binary Attribute Distribution

Distribution for attribute 'Sex':

Distribution for attribute 'ChestPainType':

Distribution for attribute 'FastingBloodSugar':

Distribution for attribute 'RestingECG':

Distribution for attribute 'ExerciseAngina':

Distribution for attribute 'ST_Slope':

Distribution for attribute 'Target':

```
1 #Visualizing the overall distribution of the numeric attributes.
2 plt.figure(figsize=(15,15))
3
4 ax=1
5 for col in ('Age', 'RestingBP_s', 'Cholesterol', 'MaxHeartRate', 'OldPeak'):
6     plt.subplot(3,2,ax)
7     plt.title(col)
8     sns.kdeplot(x=df[col])
```

```
 9      ax += 1
10
11 plt.suptitle('Numeric Attribute Distribution',y=1.01, size = 16, color = 'black', weight='
12 plt.tight_layout()
13 plt.savefig("numeric_dist.pdf",dpi=1200, bbox_inches='tight')
```
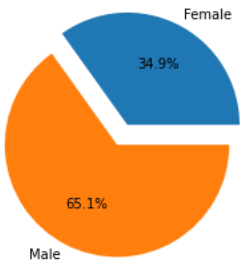
## Numeric Attribute Distribution



```
 1 #Visualizing the distribution of the categorical attributes, by target
 2 ax=1
 3 plt.figure(figsize=(15,20))
 4 plt.axis('equal')
 5
 6 for col in ('Sex','ChestPainType', 'FastingBloodSugar', 'RestingECG', 'ExerciseAngina','ST
 7   plt.subplot(6,2,ax)
 8   these_labels = labels[col]
 9   plt.title(f"Distribution for attribute '{col}' Without Heart Disease:")
10   plt.pie(no_heart_disease[col].value_counts().sort_index(),
11         autopct = '%1.1f%%', labels=these_labels,
12         explode=tuple([0.1] * len(these_labels)))
13   ax+=1
14   plt.subplot(6,2,ax)
15   these_labels = labels[col]
16   plt.title(f"Distribution for attribute '{col}' With Heart Disease:")
17   plt.pie(with_heart_disease[col].value_counts().sort_index(),
18         autopct = '%1.1f%%', labels=these_labels,
19         explode=tuple([0.1] * len(these_labels)))
20   ax+=1
21
22 plt.suptitle('Nominal/Binary Attribute Distribution by Target',y=1.01, size = 16, color =
23 plt.tight_layout()
24 plt.savefig("nominal_dist_by_target.pdf",dpi=1200, bbox_inches='tight')
```
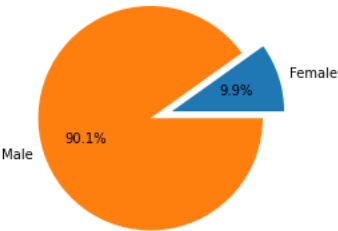
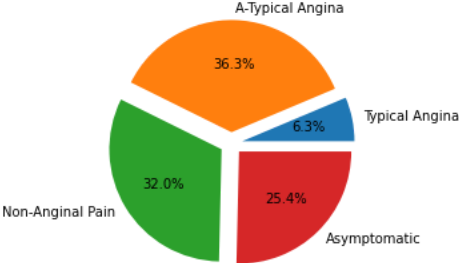## Nominal/Binary Attribute Distribution by Target

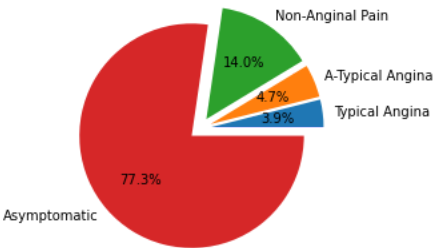Distribution for attribute 'Sex' Without Heart Disease:

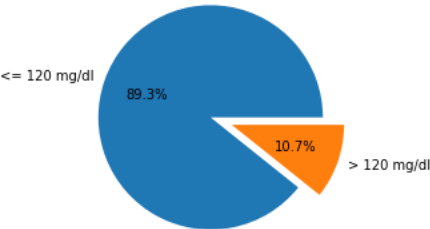Distribution for attribute 'Sex' With Heart Disease:



Distribution for attribute 'ChestPainType' Without Heart Disease:
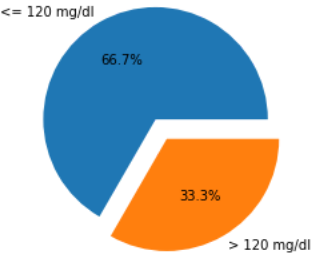
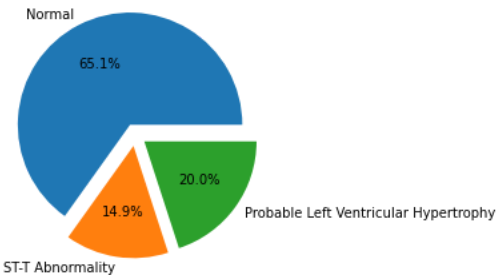Distribution for attribute 'ChestPainType' With Heart Disease:



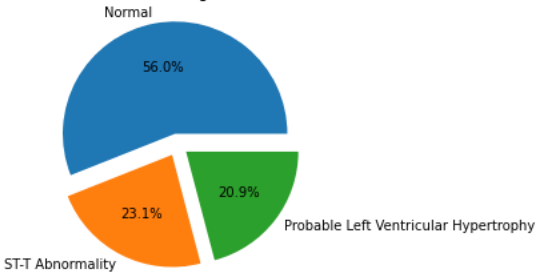Distribution for attribute 'FastingBloodSugar' Without Heart Disease:

Distribution for attribute 'FastingBloodSugar' With Heart Disease:
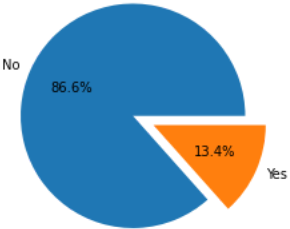


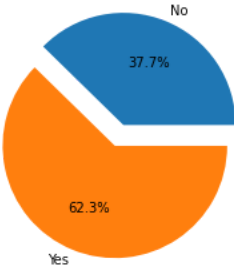Distribution for attribute 'RestingECG' Without Heart Disease:

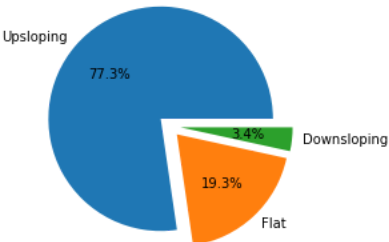Distribution for attribute 'RestingECG' With Heart Disease:



Distribution for attribute 'ExerciseAngina' Without Heart Disease:

Distribution for attribute 'ExerciseAngina' With Heart Disease:



Distribution for attribute 'ST_Slope' Without Heart Disease:

Distribution for attribute 'ST_Slope' With Heart Disease:

```python
1  #Visualizing the distribution of the numeric attributes by Target:
2  plt.figure(figsize=(15,15))
3
4  ax=1
5  for col in ('Age', 'RestingBP_s', 'Cholesterol', 'MaxHeartRate', 'OldPeak'):
6      plt.subplot(3,2,ax)
7      plt.title(col)
8      sns.kdeplot(x=no_heart_disease[col],label = "No Heart Disease")
9      sns.kdeplot(x=with_heart_disease[col],label = "Has Heart Disease")
10     plt.legend()
11     ax += 1
12
13 plt.suptitle('Numeric Attribute Distribution by Target',y=1.01, size = 16, color = 'black'
14 plt.tight_layout()
15 plt.savefig("numeric_dist_by_target.pdf",dpi=1200, bbox_inches='tight')
```

**Numeric Attribute Distribution by Target**



```
1 #Visualizing the distribution of the categorical attributes, by target and by sex
2 #For report purposes, breaking this up into two separate pages.
3 #This one for attributes 'ChestPainType', 'FastingBloodSugar', 'RestingECG'
4 #
5 ax=1
6 plt.figure(figsize=(15,30))
7 plt.axis('equal')
8
9 for col in ('ChestPainType', 'FastingBloodSugar', 'RestingECG'):
10    plt.subplot(12,2,ax)
11    these_labels = labels[col]
12    plt.title(f"Distribution for attribute '{col}'\nWithout Heart Disease(Males):")
13    plt.pie(no_heart_disease_male[col].value_counts().sort_index(),
14        autopct = '%1.1f%%', labels=these_labels,
```

```
15          explode=tuple([0.1] * len(these_labels)))
16    ax+=1
17
18    plt.subplot(12,2,ax)
19    these_labels = labels[col]
20    plt.title(f"Distribution for attribute '{col}'\nWithout Heart Disease(Females):")
21    plt.pie(no_heart_disease_female[col].value_counts().sort_index(),
22          autopct = '%1.1f%%', labels=these_labels,
23          explode=tuple([0.1] * len(these_labels)))
24    ax+=1
25
26    plt.subplot(12,2,ax)
27    these_labels = labels[col]
28    plt.title(f"Distribution for attribute '{col}'\nWith Heart Disease(Males):")
29    plt.pie(with_heart_disease_male[col].value_counts().sort_index(),
30          autopct = '%1.1f%%', labels=these_labels,
31          explode=tuple([0.1] * len(these_labels)))
32    ax+=1
33
34    plt.subplot(12,2,ax)
35    these_labels = labels[col]
36    plt.title(f"Distribution for attribute '{col}'\nWith Heart Disease(Females):")
37    plt.pie(with_heart_disease_female[col].value_counts().sort_index(),
38          autopct = '%1.1f%%', labels=these_labels,
39          explode=tuple([0.1] * len(these_labels)))
40    ax+=1
41
42  plt.suptitle('Nominal/Binary Attribute Distribution by Target and by Sex',y=1.01, size = 1
43  plt.tight_layout()
44  plt.savefig("nominal_dist_by_target_by_sex1.pdf",dpi=1200, bbox_inches='tight')
```

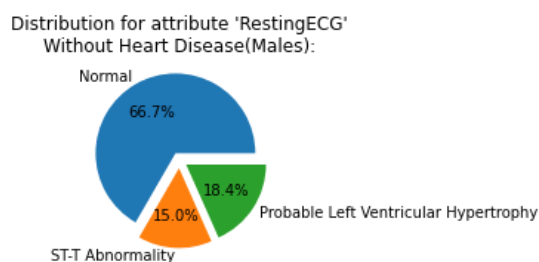## Nominal/Binary Attribute Distribution by Target and by Sex



```
 1  #Visualizing the distribution of the categorical attributes, by target and by sex
 2  #For report purposes, breaking this up into two separate pages.
 3  #This one for attributes 'ExerciseAngina','ST_Slope'
 4  #
 5  ax=1
 6  plt.figure(figsize=(15,30))
 7  plt.axis('equal')
 8
 9  for col in ('ExerciseAngina','ST_Slope'):
10      plt.subplot(12,2,ax)
11      these_labels = labels[col]
12      plt.title(f"Distribution for attribute '{col}'\nWithout Heart Disease(Males):")
```

```python
13    plt.pie(no_heart_disease_male[col].value_counts().sort_index(),
14        autopct = '%1.1f%%', labels=these_labels,
15        explode=tuple([0.1] * len(these_labels)))
16    ax+=1
17
18    plt.subplot(12,2,ax)
19    these_labels = labels[col]
20    plt.title(f"Distribution for attribute '{col}'\nWithout Heart Disease(Females):")
21    plt.pie(no_heart_disease_female[col].value_counts().sort_index(),
22        autopct = '%1.1f%%', labels=these_labels,
23        explode=tuple([0.1] * len(these_labels)))
24    ax+=1
25
26    plt.subplot(12,2,ax)
27    these_labels = labels[col]
28    plt.title(f"Distribution for attribute '{col}'\nWith Heart Disease(Males):")
29    plt.pie(with_heart_disease_male[col].value_counts().sort_index(),
30        autopct = '%1.1f%%', labels=these_labels,
31        explode=tuple([0.1] * len(these_labels)))
32    ax+=1
33
34    plt.subplot(12,2,ax)
35    these_labels = labels[col]
36    plt.title(f"Distribution for attribute '{col}'\nWith Heart Disease(Females):")
37    plt.pie(with_heart_disease_female[col].value_counts().sort_index(),
38        autopct = '%1.1f%%', labels=these_labels,
39        explode=tuple([0.1] * len(these_labels)))
40    ax+=1
41
42 plt.suptitle('Nominal/Binary Attribute Distribution by Target and by Sex, Cont\'d',y=1.01,
43 plt.tight_layout()
44 plt.savefig("nominal_dist_by_target_by_sex2.pdf",dpi=1200, bbox_inches='tight')
```

## Nominal/Binary Attribute Distribution by Target and by Sex, Cont'd

Distribution for attribute 'ExerciseAngina'
Without Heart Disease(Males):

Distribution for attribute 'ExerciseAngina'
Without Heart Disease(Females):

No 85.4%
14.6% Yes

No 88.8%
11.2% Yes

Distribution for attribute 'ExerciseAngina'
With Heart Disease(Males):

Distribution for attribute 'ExerciseAngina'
With Heart Disease(Females):

No 36.8%
63.2% Yes

No 46.0%
54.0% Yes

Distribution for attribute 'ST_Slope'
Without Heart Disease(Males):

Distribution for attribute 'ST_Slope'
Without Heart Disease(Females):

Upsloping 79.4%
4.5% Downsloping

Upsloping 73.4%
1.4% Downsloping

```
1  #Visualizing the distribution of the numerical attributes, by target and by sex
2  plt.figure(figsize=(15,15))
3
4  ax=1
5  for col in ('Age', 'RestingBP_s', 'Cholesterol', 'MaxHeartRate', 'OldPeak'):
6      plt.subplot(3,2,ax)
7      plt.title(col)
8      sns.kdeplot(x=no_heart_disease_male[col],label = "No Heart Disease(Male)")
9      sns.kdeplot(x=with_heart_disease_male[col],label = "Has Heart Disease(Male)")
10     sns.kdeplot(x=no_heart_disease_female[col],label = "No Heart Disease(Female)")
11     sns.kdeplot(x=with_heart_disease_female[col],label = "Has Heart Disease(Female)")
12     plt.legend()
13     ax += 1
14
15 plt.suptitle('Numeric Attribute Distribution, by Target and Sex',y=1.01, size = 16, color
16 plt.tight_layout()
17 plt.savefig("numeric_dist_by_target_by_sex.pdf",dpi=1200, bbox_inches='tight')
```

**Numeric Attribute Distribution, by Target and Sex**



## Outlier Detection:

```
1 #Check for outliers on numeric attributes
2 #Using for outlier detection three methods.
3 #Note: for the next stage in this project, Module 3,
4 #one or more of these outlier detection methods will be used.
5 #For now, we only want to see how many outliers per attribute are detected with each appro
6 #
7 #Methods:
8 #  #1 1.5IQR range
9 #  #2 mean +/- 3*ST-DEV (same as GT Absolute(Z-Score))
10 #  #3 Rejecting those with a value of zero (based on visualization, only needed for 'Chole
11
12 def IQR1_5_upper(data, col):
13   Q3 = np.quantile(data[col], 0.75)
14   Q1 = np.quantile(data[col], 0.25)
15   IQR = Q3 - Q1
16   return(Q3+(1.5*IQR))
17
18 def IQR1_5_lower(data, col):
19   Q3 = np.quantile(data[col], 0.75)
20   Q1 = np.quantile(data[col], 0.25)
21   IQR = Q3 - Q1
22   return(Q1-(1.5*IQR))
23
24
25 for col in ('Age', 'RestingBP_s', 'Cholesterol', 'MaxHeartRate', 'OldPeak'):
26   upper1 = IQR1_5_upper(df,col)
27   lower1 = IQR1_5_lower(df,col)
28   stdev3 = 3*df[col].std()
29   mean = df[col].mean()
30   upper2 = mean + stdev3
31   lower2 = mean - stdev3
32
33   these_outliers1 = df[(df[col] < lower1) | (df[col] > upper1)]
34   these_outliers2 = df[(df[col] < lower2) | (df[col] > upper2)]
35   these_outliers3 = df[df[col] == 0]
36
37   print(f"For attribute '{col}': The mean is {mean}, stdev3 is {stdev3}")
38   print(f"For 1.5IQR the lower range is {lower1} the upper range is {upper1}")
39   print(f"For mean +/- 3STDEV the lower range is {lower2} the upper range is {upper2}")
40   print(f"\n")
41
42   print(f"Using 1.5IQR Method:")
43   if (these_outliers1.shape[0] > 1):
44     print(f"For attribute '{col}': There are {these_outliers1.shape[0]} outliers:\n")
45     print(these_outliers1)
46     print("\n")
47   elif (these_outliers1.shape[0] == 1):
48     print(f"For attribute '{col}': There is {these_outliers1.shape[0]} outlier:\n")
49     print(these_outliers1)
50     print("\n")
51   else:
```

```
52      print(f"For attribute '{col}': There are no outliers.\n")
53      print("\n")
54
55    print(f"Using mean +/- 3STDEV Method:")
56    if (these_outliers2.shape[0] > 1):
57      print(f"For attribute '{col}': There are {these_outliers2.shape[0]} outliers:\n")
58      print(these_outliers2)
59      print("\n")
60    elif (these_outliers2.shape[0] == 1):
61      print(f"For attribute '{col}': There is {these_outliers2.shape[0]} outlier:\n")
62      print(these_outliers2)
63      print("\n")
64    else:
65      print(f"For attribute '{col}': There are no outliers.\n")
66      print("\n")
67
68    if(col == 'Cholesterol'):
69      print(f"Identifying 'zero' values(for 'Cholesterol') Method:")
70      if (these_outliers3.shape[0] > 1):
71        print(f"For attribute '{col}': There are {these_outliers3.shape[0]} outliers:\n")
72        print(these_outliers3)
73        print("\n")
74      elif (these_outliers3.shape[0] == 1):
75        print(f"For attribute '{col}': There is {these_outliers3.shape[0]} outlier:\n")
76        print(these_outliers3)
77        print("\n")
78      else:
79        print(f"For attribute '{col}': There are no outliers.\n")
80        print("\n")
81
82 print("\n\n")
```

```
For attribute 'Age': The mean is 53.49509269356598, stdev3 is 28.276802628148687
For 1.5IQR the lower range is 27.5 the upper range is 79.5
For mean +/- 3STDEV the lower range is 25.218290065417293 the upper range is 81.77189


Using 1.5IQR Method:
For attribute 'Age': There are no outliers.



Using mean +/- 3STDEV Method:
For attribute 'Age': There are no outliers.



For attribute 'RestingBP_s': The mean is 132.3773173391494, stdev3 is 55.545341581484
For 1.5IQR the lower range is 90.0 the upper range is 170.0
For mean +/- 3STDEV the lower range is 76.83197575766476 the upper range is 187.92265


Using 1.5IQR Method:
```

```
For attribute 'RestingBP_s': There are 28 outliers:
```

|     | Age | Sex | ChestPainType | RestingBP_s | Cholesterol | FastingBloodSugar \ |
|-----|-----|-----|---------------|-------------|-------------|---------------------|
| 109 | 39  | 1   | 2             | 190         | 241         | 0                   |
| 123 | 58  | 0   | 2             | 180         | 393         | 0                   |
| 189 | 53  | 1   | 4             | 180         | 285         | 0                   |
| 190 | 46  | 1   | 4             | 180         | 280         | 0                   |
| 241 | 54  | 1   | 4             | 200         | 198         | 0                   |
| 274 | 45  | 0   | 2             | 180         | 295         | 0                   |
| 275 | 59  | 1   | 3             | 180         | 213         | 0                   |
| 278 | 57  | 0   | 4             | 180         | 347         | 0                   |
| 314 | 53  | 1   | 4             | 80          | 0           | 0                   |
| 365 | 64  | 0   | 4             | 200         | 0           | 0                   |
| 372 | 63  | 1   | 4             | 185         | 0           | 0                   |
| 399 | 61  | 1   | 3             | 200         | 0           | 1                   |
| 411 | 54  | 1   | 4             | 180         | 0           | 1                   |
| 423 | 60  | 1   | 3             | 180         | 0           | 0                   |
| 449 | 55  | 1   | 3             | 0           | 0           | 0                   |
| 475 | 59  | 1   | 4             | 178         | 0           | 1                   |
| 550 | 55  | 1   | 4             | 172         | 260         | 0                   |
| 585 | 57  | 1   | 2             | 180         | 285         | 1                   |
| 592 | 61  | 1   | 4             | 190         | 287         | 1                   |
| 673 | 59  | 0   | 4             | 174         | 249         | 0                   |
| 702 | 59  | 1   | 1             | 178         | 270         | 0                   |
| 725 | 55  | 0   | 4             | 180         | 327         | 0                   |
| 732 | 56  | 0   | 4             | 200         | 288         | 1                   |
| 759 | 54  | 1   | 2             | 192         | 283         | 0                   |
| 774 | 66  | 0   | 4             | 178         | 228         | 1                   |
| 780 | 64  | 0   | 4             | 180         | 325         | 0                   |
| 855 | 68  | 1   | 3             | 180         | 274         | 1                   |
| 880 | 52  | 1   | 3             | 172         | 199         | 1                   |

|     | RestingECG | MaxHeartRate | ExerciseAngina | OldPeak | ST_Slope | Target |
|-----|------------|--------------|----------------|---------|----------|--------|
| 109 | 0          | 106          | 0              | 0.0     | 1        | 0      |
| 123 | 0          | 110          | 1              | 1.0     | 2        | 1      |
| 189 | 1          | 120          | 1              | 1.5     | 2        | 1      |

## Data Manipulation, for Outliers, and Model Considerations:

```
 1 #Before we can run multinomial Naive Bayes we must remove any negative numbers in the data
 2 mins = df.min()
 3 print(mins)
 4 print("\n\n")
 5 #There are only negative values for attribute 'OldPeak'
 6 #Applying a simple shift to eliminate any negatives.
 7
 8 df_no_negs = df.copy()
 9 df_no_negs.name = "No Negatives Data-Set"
10
11 df_no_negs['OldPeak'] = df_no_negs['OldPeak'] + abs(df_no_negs['OldPeak'].min())
12
```

```
13 #Visualizing the overall distribution of 'OldPeak' before and after modification for negat
14 plt.figure(figsize=(10,5))
15
16 plt.subplot(1,2,1)
17 plt.title('OldPeak')
18 sns.kdeplot(x=df['OldPeak'])
19
20 plt.subplot(1,2,2)
21 plt.title('OldPeakNoNeg')
22 sns.kdeplot(x=df_no_negs['OldPeak'])
23
24 plt.suptitle('Visualizing the Transformation for OldPeak',y=1.01, size = 16, color = 'blac
25 plt.tight_layout()
26 plt.savefig("oldpeak_transformation.pdf",dpi=1200, bbox_inches='tight')
27
```

```
Age                28.0
Sex                 0.0
ChestPainType       1.0
RestingBP_s         0.0
Cholesterol         0.0
FastingBloodSugar   0.0
RestingECG          0.0
MaxHeartRate       60.0
ExerciseAngina      0.0
OldPeak            -2.6
ST_Slope            1.0
Target              0.0
dtype: float64
```



Visualizing the Transformation for OldPeak

```
 1 #Only addressing outliers for attribute Cholesterol, specifically the instances where Chol
 2 #All other outliers appear in much smaller quantities.
 3
 4 df_outliers_addressed = df_no_negs.copy()
 5 df_outliers_addressed.name = "Outliers Addressed Data-Set"
 6
 7 cholesterol_mean = df_outliers_addressed['Cholesterol'][df_outliers_addressed['Cholesterol
 8
 9 df_outliers_addressed['Cholesterol'].replace(to_replace=0.0, value=cholesterol_mean, inpla
10
11
12 #Visualizing the overall distribution Cholesterol before and after dealing with outliers.
13 plt.figure(figsize=(10,5))
14
15 plt.subplot(1,2,1)
16 plt.title('Cholesterol')
17 sns.kdeplot(x=df['Cholesterol'])
18
19 plt.subplot(1,2,2)
20 plt.title('Cholesterol(No Outliers)')
21 sns.kdeplot(x=df_outliers_addressed['Cholesterol'])
22
23 plt.suptitle('Visualizing the Transformation for Cholesterol',y=1.01, size = 16, color = '
24 plt.tight_layout()
```



**Visualizing the Transformation for Cholesterol**

### Normalize Numeric Data for potential model training:

```
 1 #Normalizing numeric data to see if this helps, or hinders the accuracy of the ML models.
 2 df_normalized = df_outliers_addressed.copy()
 3 df_normalized.name = "Normalized Data-Set"
```

```
  4
  5 for col in ('Age', 'RestingBP_s', 'Cholesterol', 'MaxHeartRate', 'OldPeak'):
  6   df_normalized[col] = (df[col]-df[col].min())/(df[col].max()-df[col].min())
```

## Create ONE-HOT columns for all categorical attributes:

```
 1 #For Nominal & Binary (ie categorical) attributes, perform one-hot conversion.
 2 #convert only categorical variables/features to dummy/one-hot features
 3 cat_cols = ['Sex','ChestPainType', 'FastingBloodSugar', 'RestingECG', 'ExerciseAngina','ST
 4
 5 df_onehot = pd.get_dummies(df, columns=cat_cols, prefix = cat_cols)
 6 df_onehot.name = "Original Data-Set, After ONEHOT"
 7
 8 df_onehot_no_negs = pd.get_dummies(df_no_negs, columns=cat_cols, prefix = cat_cols)
 9 df_onehot_no_negs.name = "No negatives Data-Set, After ONEHOT"
10
11 df_onehot_outliers_addressed = pd.get_dummies(df_outliers_addressed, columns=cat_cols, pre
12 df_onehot_outliers_addressed.name = "Outliers Addressed Data-Set, After ONEHOT"
13
14 df_onehot_normalized = pd.get_dummies(df_normalized, columns=cat_cols, prefix = cat_cols)
15 df_onehot_normalized.name = "Normalized Data-Set, After ONEHOT"
16
```

## Use SMOTE to Balance the Class Variable:

```
 1 # Balance the dataset using SMOTE.
 2
 3 def do_SMOTE(df, classCol="Target"):
 4   oversample = SMOTE()
 5   X = df.drop(columns=classCol)
 6   Y = df[classCol]
 7   X, Y = oversample.fit_resample(X, Y)
 8   X[classCol] = Y
 9   X.name = df.name+", After SMOTE"
10   return(X)
11
12
13
14 df_smote = do_SMOTE(df)
15 df_no_negs_smote = do_SMOTE(df_no_negs)
16
17 df_normalized_smote = do_SMOTE(df_normalized)
18 df_outliers_addressed_smote = do_SMOTE(df_outliers_addressed)
19
20 df_onehot_smote = do_SMOTE(df_onehot)
21 df_onehot_no_negs_smote = do_SMOTE(df_onehot_no_negs)
22 df_onehot_outliers_addressed_smote = do_SMOTE(df_onehot_outliers_addressed)
```

```
23 df_onehot_normalized_smote = do_SMOTE(df_onehot_normalized)
24
25
26 #Visualizing the distribution of the class variable before and after SMOTE for the basic d
27 ax=1
28 plt.figure(figsize=(10,6))
29
30 for this_df in (df,df_smote):
31   plt.subplot(1,2,ax)
32   these_labels = labels['Target']
33   plt.title(f"{this_df.name}:")
34   plt.pie(this_df['Target'].value_counts().sort_index(),
35         autopct = '%1.1f%%', labels=these_labels,
36         explode=tuple([0.1] * len(these_labels)))
37   plt.axis('equal')
38   ax+=1
39
40 plt.suptitle('Distribution for Class Variable Before and After SMOTE',y=1.01, size = 16, c
41 plt.tight_layout()
42 plt.savefig("smote.pdf",dpi=1200, bbox_inches='tight')
43
```

**Distribution for Class Variable Before and After SMOTE**

Original Data-Set: | Original Data-Set, After SMOTE:

Without Heart Disease | Without Heart Disease

44.7% | 50.0%

55.3% | 50.0%

Has Heart Disease | Has Heart Disease

## ML Algorithms:

```
1 def do_DT(df,levels,class_col_name,verbose=0):
2   #not disabling randomness.
```

```
 3    #np.random.seed(0)

 4

 5    # Split dataset into training set and test set
 6    feature_names=df.columns[df.columns != class_col_name ]
 7    # 80% training and 20% test
 8    X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c

 9

10    clf = tree.DecisionTreeClassifier(max_depth=levels,criterion='gini')
11    clf = clf.fit(X_train, Y_train)
12    if (verbose >= 1):
13      print(f"Successfuly trained the decision tree for {levels} levels...")

14

15    # Let's make the prdictions on the test set  that we set aside earlier using the trained
16    Y_pred = clf.predict(X_test)

17

18    cf=confusion_matrix(Y_test, Y_pred)
19    tn, fp, fn, tp=cf.ravel()
20    tpr=0.0
21    fpr=0.0
22    tpr = tp/(tp+fp)
23    fpr = fp/(fp+tn)
24    fnr = fn/(fn+tp)

25

26    if (verbose >= 2):
27      print ("Confusion Matrix")
28      print(cf)
29      print("")
30      print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
31      print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)

32

33    #print precision, recall, and accuracy from the perspective of each of the class (0 and
34    if (verbose >= 2):
35      print(classification_report(Y_test, Y_pred, digits=3))

36

37    accuracy = accuracy_score(Y_test, Y_pred)
38    f1_weighted = f1_score(Y_test, Y_pred,average='weighted')

39

40    if (verbose >= 1):
41      print(f"Accuracy is: {accuracy}")
42      print(f"F1 Weighted is: {f1_weighted}")
43      print("")

44

45    return(accuracy,f1_weighted,tpr,fpr,fnr)


 1 def do_mnNB(df,class_col_name,verbose=0):
 2    #not disabling randomness.
 3    #np.random.seed(0)

 4

 5    # Split dataset into training set and test set
 6    feature_names=df.columns[df.columns != class_col_name ]
```

```
 7    # 80% training and 20% test
 8    X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c

 9
10    #Create a MultiNomial NB Classifier
11    nb = MultinomialNB()

12
13    #Train the model using the training sets
14    nb.fit(X_train, Y_train)

15
16    #Predict the response for test dataset
17    Y_pred = nb.predict(X_test)

18
19    if (verbose >= 2):
20      print ("Total Columns (including class)",len(df.columns))
21      print("Classes ",nb.classes_)
22      print("Number of records for classes ",nb.class_count_)
23      print("Log prior probability for classes ", nb.class_log_prior_)
24      print("Log conditional probability for each feature given a class\n",nb.feature_log_pr

25
26    cf=confusion_matrix(Y_test, Y_pred)
27    tn, fp, fn, tp=cf.ravel()
28    tpr = tp/(tp+fp)
29    fpr = fp/(fp+tn)
30    fnr = fn/(fn+tp)

31
32    if (verbose >= 2):
33      print ("Confusion Matrix")
34      print(cf)
35      print("")
36      print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
37      print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)

38
39    if (verbose >= 2):
40      print(classification_report(Y_test, Y_pred, digits=3))

41
42    accuracy = accuracy_score(Y_test, Y_pred)
43    f1_weighted = f1_score(Y_test, Y_pred,average='weighted')

44
45    if (verbose >= 1):
46      print(f"Accuracy is: {accuracy}")
47      print(f"F1 Weighted is: {f1_weighted}")
48      print("")

49
50    return(accuracy,f1_weighted,tpr,fpr,fnr)



 1 def do_gaNB(df,class_col_name,verbose=0):
 2   #not disabling randomness.
 3   #np.random.seed(0)
 4
 5   # Split dataset into training set and test set
```

```
 6    feature_names=df.columns[df.columns != class_col_name ]
 7    # 80% training and 20% test
 8    X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c
 9
10    #Create a Gaussian NB Classifier
11    nb = GaussianNB()
12
13    #Train the model using the training sets
14    nb.fit(X_train, Y_train)
15
16    #Predict the response for test dataset
17    Y_pred = nb.predict(X_test)
18
19    if (verbose >= 2):
20      print ("Total Columns (including class)",len(df.columns))
21      print("Number of records for classes ",nb.class_count_)
22
23    cf=confusion_matrix(Y_test, Y_pred)
24    tn, fp, fn, tp=cf.ravel()
25    tpr = tp/(tp+fp)
26    fpr = fp/(fp+tn)
27    fnr = fn/(fn+tp)
28
29    if (verbose >= 2):
30      print ("Confusion Matrix")
31      print(cf)
32      print("")
33      print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
34      print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)
35
36    if (verbose >= 2):
37      print(classification_report(Y_test, Y_pred, digits=3))
38
39    accuracy = accuracy_score(Y_test, Y_pred)
40    f1_weighted = f1_score(Y_test, Y_pred,average='weighted')
41
42    if (verbose >= 1):
43      print(f"Accuracy is: {accuracy}")
44      print(f"F1 Weighted is: {f1_weighted}")
45      print("")
46
47    return(accuracy,f1_weighted,tpr,fpr,fnr)



 1 def do_LR(df,class_col_name,verbose=0):
 2   #not disabling randomness.
 3   #np.random.seed(0)
 4
 5   # Split dataset into training set and test set
 6   feature_names=df.columns[df.columns != class_col_name ]
 7   # 80% training and 20% test
```

```
8    X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c
9
10   lr = LogisticRegression(max_iter=2000)
11
12   #Train the model using the training sets
13   lr.fit(X_train, Y_train)
14
15   #Predict the response for test dataset
16   Y_pred = lr.predict(X_test)
17
18   if (verbose >= 2):
19     print ("Total Columns (including class)",len(df.columns))
20
21   cf=confusion_matrix(Y_test, Y_pred)
22   tn, fp, fn, tp=cf.ravel()
23   tpr = tp/(tp+fp)
24   fpr = fp/(fp+tn)
25   fnr = fn/(fn+tp)
26
27   if (verbose >= 2):
28     print ("Confusion Matrix")
29     print(cf)
30     print("")
31     print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
32     print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)
33
34   if (verbose >= 2):
35     print(classification_report(Y_test, Y_pred, digits=3))
36
37   accuracy = accuracy_score(Y_test, Y_pred)
38   f1_weighted = f1_score(Y_test, Y_pred,average='weighted')
39
40   if (verbose >= 1):
41     print(f"Accuracy is: {accuracy}")
42     print(f"F1 Weighted is: {f1_weighted}")
43     print("")
44
45   return(accuracy,f1_weighted,tpr,fpr,fnr)
```

```
1 def do_KNN(df,class_col_name,verbose=0):
2   #not disabling randomness.
3   #np.random.seed(0)
4
5   # Split dataset into training set and test set
6   feature_names=df.columns[df.columns != class_col_name ]
7   # 80% training and 20% test
8   X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c
9
10  knn = KNeighborsClassifier()
11
```

```
12    #Train the model using the training sets
13    knn.fit(X_train, Y_train)
14
15    #Predict the response for test dataset
16    Y_pred = knn.predict(X_test)
17
18    if (verbose >= 2):
19      print ("Total Columns (including class)",len(df.columns))
20
21    cf=confusion_matrix(Y_test, Y_pred)
22    tn, fp, fn, tp=cf.ravel()
23    tpr = tp/(tp+fp)
24    fpr = fp/(fp+tn)
25    fnr = fn/(fn+tp)
26
27    if (verbose >= 2):
28      print ("Confusion Matrix")
29      print(cf)
30      print("")
31      print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
32      print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)
33
34    if (verbose >= 2):
35      print(classification_report(Y_test, Y_pred, digits=3))
36
37    accuracy = accuracy_score(Y_test, Y_pred)
38    f1_weighted = f1_score(Y_test, Y_pred,average='weighted')
39
40    if (verbose >= 1):
41      print(f"Accuracy is: {accuracy}")
42      print(f"F1 Weighted is: {f1_weighted}")
43      print("")
44
45    return(accuracy,f1_weighted,tpr,fpr,fnr)
```

```
1 def do_RF(df,class_col_name,verbose=0):
2    #not disabling randomness.
3    #np.random.seed(0)
4
5    # Split dataset into training set and test set
6    feature_names=df.columns[df.columns != class_col_name ]
7    # 80% training and 20% test
8    X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c
9
10   rf = RandomForestClassifier()
11
12   #Train the model using the training sets
13   rf.fit(X_train, Y_train)
14
15   #Predict the response for test dataset
```

```
16   Y_pred = rf.predict(X_test)

17

18   if (verbose >= 2):
19     print ("Total Columns (including class)",len(df.columns))

20

21   cf=confusion_matrix(Y_test, Y_pred)
22   tn, fp, fn, tp=cf.ravel()
23   tpr = tp/(tp+fp)
24   fpr = fp/(fp+tn)
25   fnr = fn/(fn+tp)

26

27   if (verbose >= 2):
28     print ("Confusion Matrix")
29     print(cf)
30     print("")
31     print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
32     print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)

33

34   if (verbose >= 2):
35     print(classification_report(Y_test, Y_pred, digits=3))

36

37   accuracy = accuracy_score(Y_test, Y_pred)
38   f1_weighted = f1_score(Y_test, Y_pred,average='weighted')

39

40   if (verbose >= 1):
41     print(f"Accuracy is: {accuracy}")
42     print(f"F1 Weighted is: {f1_weighted}")
43     print("")

44

45   return(accuracy,f1_weighted,tpr,fpr,fnr)
```

```
 1 def do_SVM(df,class_col_name,verbose=0):
 2   #not disabling randomness.
 3   #np.random.seed(0)

 4

 5   # Split dataset into training set and test set
 6   feature_names=df.columns[df.columns != class_col_name ]
 7   # 80% training and 20% test
 8   X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c

 9

10   svm = SVC()

11

12   #Train the model using the training sets
13   svm.fit(X_train, Y_train)

14

15   #Predict the response for test dataset
16   Y_pred = svm.predict(X_test)

17

18   if (verbose >= 2):
19     print ("Total Columns (including class)",len(df.columns))
```

```
20
21   cf=confusion_matrix(Y_test, Y_pred)
22   tn, fp, fn, tp=cf.ravel()
23   tpr = tp/(tp+fp)
24   fpr = fp/(fp+tn)
25   fnr = fn/(fn+tp)
26
27
28   if (verbose >= 2):
29     print ("Confusion Matrix")
30     print(cf)
31     print("")
32     print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
33     print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)
34
35   if (verbose >= 2):
36     print(classification_report(Y_test, Y_pred, digits=3))
37
38   accuracy = accuracy_score(Y_test, Y_pred)
39   f1_weighted = f1_score(Y_test, Y_pred,average='weighted')
40
41   if (verbose >= 1):
42     print(f"Accuracy is: {accuracy}")
43     print(f"F1 Weighted is: {f1_weighted}")
44     print("")
45
46   return(accuracy,f1_weighted,tpr,fpr,fnr)
```

```
 1 def do_XGB(df,class_col_name,verbose=0):
 2   #not disabling randomness.
 3   #np.random.seed(0)
 4
 5   # Split dataset into training set and test set
 6   feature_names=df.columns[df.columns != class_col_name ]
 7   # 80% training and 20% test
 8   X_train, X_test, Y_train, Y_test = train_test_split(df.loc[:, feature_names], df[class_c
 9
10   xgb = XGBClassifier()
11
12   #Train the model using the training sets
13   xgb.fit(X_train, Y_train)
14
15   #Predict the response for test dataset
16   Y_pred = xgb.predict(X_test)
17
18   if (verbose >= 2):
19     print ("Total Columns (including class)",len(df.columns))
20
21   cf=confusion_matrix(Y_test, Y_pred)
22   tn, fp, fn, tp=cf.ravel()
```

```
23   tpr = tp/(tp+fp)
24   fpr = fp/(fp+tn)
25   fnr = fn/(fn+tp)
26
27   if (verbose >= 2):
28     print ("Confusion Matrix")
29     print(cf)
30     print("")
31     print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
32     print ("TPR: ",tpr,", FPR: ",fpr, "FNR: ",fnr)
33
34   if (verbose >= 2):
35     print(classification_report(Y_test, Y_pred, digits=3))
36
37   accuracy = accuracy_score(Y_test, Y_pred)
38   f1_weighted = f1_score(Y_test, Y_pred,average='weighted')
39
40   if (verbose >= 1):
41     print(f"Accuracy is: {accuracy}")
42     print(f"F1 Weighted is: {f1_weighted}")
43     print("")
44
45   return(accuracy,f1_weighted,tpr,fpr,fnr)
```

## Initial Run of All ML Algorithms:

```
1  #Initial Run of all ML Algorithms just to make sure everything works correctly.
2  #Using original data:
3
4  for i in range(3,11):
5    print(f"DT with {i} levels:")
6    do_DT(df,i,'Target',5)
7
8  print(f"MN_NB:")
9  do_mnNB(df_no_negs,'Target',5)
10
11 print(f"GA_NB:")
12 do_gaNB(df,'Target',5)
13
14 print(f"LR:")
15 do_LR(df,'Target',5)
16
17 print(f"KNN:")
18 do_KNN(df,'Target',5)
19
20 print(f"RF:")
21 do_RF(df,'Target',5)
22
23 print(f"SVM:")
```

```
24 do_SVM(df,'Target',5)
25
26 print(f"XGB:")
27 do_XGB(df,'Target',5)
```

```
    DT with 3 levels:
    Successfuly trained the decision tree for 3 levels...
    Confusion Matrix
    [[76 15]
     [13 80]]

    TP:  80 , FP:  15 , TN:  76 , FN: 13
    TPR:  0.8421052631578947 , FPR:  0.16483516483516483 FNR:  0.13978494623655913
                  precision    recall  f1-score   support

               0      0.854     0.835     0.844        91
               1      0.842     0.860     0.851        93

        accuracy                          0.848       184
       macro avg      0.848     0.848     0.848       184
    weighted avg      0.848     0.848     0.848       184

    Accuracy is: 0.8478260869565217
    F1 Weighted is: 0.8477901120361805

    DT with 4 levels:
    Successfuly trained the decision tree for 4 levels...
    Confusion Matrix
    [[71 13]
     [11 89]]

    TP:  89 , FP:  13 , TN:  71 , FN: 11
    TPR:  0.8725490196078431 , FPR:  0.15476190476190477 FNR:  0.11
                  precision    recall  f1-score   support

               0      0.866     0.845     0.855        84
               1      0.873     0.890     0.881       100

        accuracy                          0.870       184
       macro avg      0.869     0.868     0.868       184
    weighted avg      0.869     0.870     0.869       184

    Accuracy is: 0.8695652173913043
    F1 Weighted is: 0.8694251824344299

    DT with 5 levels:
    Successfuly trained the decision tree for 5 levels...
    Confusion Matrix
    [[65 16]
     [15 88]]

    TP:  88 , FP:  16 , TN:  65 , FN: 15
    TPR:  0.8461538461538461 , FPR:  0.19753086419753085 FNR:  0.14563106796116504
                  precision    recall  f1-score   support

               0      0.812     0.802     0.807        81
```

```
        1        0.846       0.854       0.850        103

   accuracy                              0.832        184
  macro avg        0.829       0.828     0.829        184
weighted avg       0.831       0.832     0.831        184


Accuracy is: 0.8315217391304348
```

## Validation:

```python
1 #Create a subroutine to invoke cross_validate for a given model, and dataset.
2 #
3 def doCV(test_num, test_name, model, this_df, folds=5, classCol='Target', verbose=0):
4   X = this_df.drop(classCol,axis=1)
5   Y = this_df[classCol]
6
7   scoring = {'acc': 'accuracy',
8              'rec': 'recall_macro'}
9   scores = cross_validate(model, X, Y, scoring=scoring,
10                       cv=folds, return_train_score=True)
11
12   ACC = scores['test_acc']
13   ACC_mean = ACC.mean() * 100
14   ACC_std = ACC.std()
15
16   TPR = scores['test_rec']
17   TPR_mean = TPR.mean() * 100
18   TPR_std = TPR.std()
19
20   FNR = 1-TPR
21   FNR_mean = FNR.mean() * 100
22   FNR_std = FNR.std()
23
24   if verbose > 0:
25     print(f"For TEST #{test_num}: {test_name:35s} After Cross-Val, Using Data-Set: {this_d
26
27   #Store all results in an easy lookup-table
28   results[f"LAST_TEST"] = test_num
29   results[f"TEST{test_num}#NAME"] = test_name
30   results[f"TEST{test_num}#DFNAME"] = this_df.name
31   results[f"TEST{test_num}#DF_ROWS"] = this_df.shape[0]
32   results[f"TEST{test_num}#ACCURACY_MEAN"] = ACC_mean
33   results[f"TEST{test_num}#ACCURACY_STD"] = ACC_std
34   results[f"TEST{test_num}#ACCURACY"] = ACC
35   results[f"TEST{test_num}#TPR_MEAN"] = TPR_mean
36   results[f"TEST{test_num}#TPR_STD"] = TPR_std
37   results[f"TEST{test_num}#FNR_MEAN"] = FNR_mean
38   results[f"TEST{test_num}#FNR_STD"] = FNR_std
39   results[f"TEST{test_num}#FNR"] = FNR
40
41 #Create subroutines to display test results, and obtain accuracy and FNR for a given test.
```

```python
42 def getTestAccuracy(x):
43   value = results[f"TEST{x}#ACCURACY_MEAN"]
44   return(value)
45
46 def getTestF1(x):
47   value = results[f"TEST{x}#F1_WEIGHTED_MEAN"]
48   return(value)
49
50 def getTestFNR(x):
51   value = results[f"TEST{x}#FNR_MEAN"]
52   return(value)
53
54 def getTestAccVec(x):
55   value = results[f"TEST{x}#ACCURACY"]
56   return(value)
57
58 def getTestFNRVec(x):
59   value = results[f"TEST{x}#FNR"]
60   return(value)
61
62 def getTestDFRows(x):
63   value = results[f"TEST{x}#DF_ROWS"]
64   return(value)
65
66
67 def displayResult(x):
68   name = results[f"TEST{x}#NAME"]
69   df_name = results[f"TEST{x}#DFNAME"]
70   accuracy = results[f"TEST{x}#ACCURACY_MEAN"]
71   #f1_weighted = results[f"TEST{x}#F1_WEIGHTED_MEAN"]
72   tpr = results[f"TEST{x}#TPR_MEAN"]
73   #fpr = results[f"TEST{x}#FPR_MEAN"]
74   fnr = results[f"TEST{x}#FNR_MEAN"]
75
76   print(f"Using ML Model: {name} and {df_name}:")
77   #print(f"TEST #{x}: Average Accuracy is {accuracy:.2f}%, Average F1(weighted) is {f1_wei
78   print(f"TEST #{x}: Average Accuracy is {accuracy:.2f}%, Average TPR is {tpr:.2f}%, Avera
79   print("")
80
81
```

```python
1 #Initial Models used:
2 dt3 = tree.DecisionTreeClassifier(max_depth=3,criterion='gini')
3 dt4 = tree.DecisionTreeClassifier(max_depth=4,criterion='gini')
4 dt5 = tree.DecisionTreeClassifier(max_depth=5,criterion='gini')
5 dt6 = tree.DecisionTreeClassifier(max_depth=6,criterion='gini')
6 dt7 = tree.DecisionTreeClassifier(max_depth=7,criterion='gini')
7 dt8 = tree.DecisionTreeClassifier(max_depth=8,criterion='gini')
8 dt9 = tree.DecisionTreeClassifier(max_depth=9,criterion='gini')
9 dt10 = tree.DecisionTreeClassifier(max_depth=10,criterion='gini')
```

```python
10 dt11 = tree.DecisionTreeClassifier(max_depth=11,criterion='gini')
11 lr = LogisticRegression(max_iter=2000)
12 knn = KNeighborsClassifier()
13 svm = SVC()
14 mn_nb = MultinomialNB()
15 ga_nb = GaussianNB()
16 rf = RandomForestClassifier()
17 xgb = XGBClassifier()
```

## Initial Validation:

```python
 1 results = {}
 2 test_num = 0
 3 folds = 10
 4 classCol = "Target"
 5 df_list = [df,df_normalized,df_onehot,df_outliers_addressed,df_smote,df_normalized_smote,d
 6
 7 for test in ["DT_3","DT_4","DT_5","DT_6","DT_7","DT_8","DT_9","DT_10","DT_11","GA_NB","LR"
 8   if test == "GA_NB":
 9     model = ga_nb
10     test_name = "Gaussian Naive Bayes(GA-NB)"
11   elif test == "LR":
12     model = lr
13     test_name = "Logistic Regression(LR)"
14   elif test == "SVM":
15     model = svm
16     test_name = "Support Vector Machines(SVM)"
17   elif test == "KNN":
18     model = knn
19     test_name = "K Nearest Neighbours(KNN)"
20   elif test == "RF":
21     model = rf
22     test_name = "Random Forest(RF)"
23   elif test == "XGB":
24     model = xgb
25     test_name = "XG Boost(XGB)"
26   elif test == "DT_3":
27     model = dt3
28     test_name = "Decision Tree: 3 levels"
29   elif test == "DT_4":
30     model = dt4
31     test_name = "Decision Tree: 4 levels"
32   elif test == "DT_5":
33     model = dt5
34     test_name = "Decision Tree: 5 levels"
35   elif test == "DT_6":
36     model = dt6
37     test_name = "Decision Tree: 6 levels"
38   elif test == "DT_7":
```

```
39      model = dt7
40      test_name = "Decision Tree: 7 levels"
41    elif test == "DT_8":
42      model = dt8
43      test_name = "Decision Tree: 8 levels"
44    elif test == "DT_9":
45      model = dt9
46      test_name = "Decision Tree: 9 levels"
47    elif test == "DT_10":
48      model = dt10
49      test_name = "Decision Tree: 10 levels"
50    elif test == "DT_11":
51      model = dt11
52      test_name = "Decision Tree: 11 levels"
53
54    for idx,df_not_used in enumerate(df_list):
55      this_df = df_list[idx]
56      doCV(test_num,test_name,model,this_df,folds,classCol,1)
57      test_num += 1
```

```
For TEST #0: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Orig
For TEST #1: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Norr
For TEST #2: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Orig
For TEST #3: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Out
For TEST #4: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Orig
For TEST #5: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Norr
For TEST #6: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Orig
For TEST #7: Decision Tree: 3 levels                After Cross-Val, Using Data-Set: Out
For TEST #8: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Orig
For TEST #9: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Norr
For TEST #10: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Or
For TEST #11: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Ou
For TEST #12: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Or
For TEST #13: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: No
For TEST #14: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Or
For TEST #15: Decision Tree: 4 levels                After Cross-Val, Using Data-Set: Ou
For TEST #16: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Or
For TEST #17: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: No
For TEST #18: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Or
For TEST #19: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Ou
For TEST #20: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Or
For TEST #21: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: No
For TEST #22: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Or
For TEST #23: Decision Tree: 5 levels                After Cross-Val, Using Data-Set: Ou
For TEST #24: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Or
For TEST #25: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: No
For TEST #26: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Or
For TEST #27: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Ou
For TEST #28: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Or
For TEST #29: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: No
For TEST #30: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Or
For TEST #31: Decision Tree: 6 levels                After Cross-Val, Using Data-Set: Ou
For TEST #32: Decision Tree: 7 levels                After Cross-Val, Using Data-Set: Or
For TEST #33: Decision Tree: 7 levels                After Cross-Val, Using Data-Set: No
```

```
     For TEST #34: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: Or
     For TEST #35: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #36: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: Or
     For TEST #37: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: No
     For TEST #38: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: Or
     For TEST #39: Decision Tree: 7 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #40: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Or
     For TEST #41: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: No
     For TEST #42: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Or
     For TEST #43: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #44: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Or
     For TEST #45: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: No
     For TEST #46: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Or
     For TEST #47: Decision Tree: 8 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #48: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Or
     For TEST #49: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: No
     For TEST #50: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Or
     For TEST #51: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #52: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Or
     For TEST #53: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: No
     For TEST #54: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Or
     For TEST #55: Decision Tree: 9 levels              After Cross-Val, Using Data-Set: Ou
     For TEST #56: Decision Tree: 10 levels             After Cross-Val, Using Data-Set: Or
```

```python
 1 #View these results again, this time just the top 10 using each sort method:
 2 range_limit = min(10,results[f"LAST_TEST"]+1) #Top 10 results desired.
 3
 4 results_list = list(range(0,results[f"LAST_TEST"]+1))
 5 results_list.sort(key=getTestAccuracy, reverse=True)
 6
 7 print("Results of ML Models: (sorted by accuracy) (top 10)")
 8 print("")
 9 for i in results_list[0:range_limit]:  #Count starts at zero, so dont need to add one.
10   displayResult(i)
11 print("")
12 print("")
13
14
15 results_list = list(range(0,results[f"LAST_TEST"]+1))
16 results_list.sort(key=getTestFNR, reverse=False)
17
18 print("Results of ML Models: (sorted by False Negative Rate(FNR)) (top 10)")
19 print("")
20 for i in results_list[0:range_limit]:  #Count starts at zero, so dont need to add one.
21   displayResult(i)
22 print("")
23 print("")
```

```
    Results of ML Models: (sorted by accuracy) (top 10)

    Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
    TEST #108: Average Accuracy is 87.57%, Average TPR is 87.62%, Average FNR is 12.38%
```

```
Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
TEST #116: Average Accuracy is 86.87%, Average TPR is 86.93%, Average FNR is 13.07%

Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
TEST #110: Average Accuracy is 86.68%, Average TPR is 86.74%, Average FNR is 13.26%

Using ML Model: Random Forest(RF) and Outliers Addressed Data-Set, After ONEHOT, Afte
TEST #111: Average Accuracy is 86.48%, Average TPR is 86.53%, Average FNR is 13.47%

Using ML Model: XG Boost(XGB) and Original Data-Set, After ONEHOT, After SMOTE:
TEST #118: Average Accuracy is 86.48%, Average TPR is 86.53%, Average FNR is 13.47%

Using ML Model: Random Forest(RF) and Normalized Data-Set, After SMOTE:
TEST #109: Average Accuracy is 86.09%, Average TPR is 86.14%, Average FNR is 13.86%

Using ML Model: XG Boost(XGB) and Normalized Data-Set, After SMOTE:
TEST #117: Average Accuracy is 85.98%, Average TPR is 86.03%, Average FNR is 13.97%

Using ML Model: Random Forest(RF) and Outliers Addressed Data-Set:
TEST #107: Average Accuracy is 85.91%, Average TPR is 85.64%, Average FNR is 14.36%

Using ML Model: Random Forest(RF) and Normalized Data-Set:
TEST #105: Average Accuracy is 85.91%, Average TPR is 85.54%, Average FNR is 14.46%

Using ML Model: Logistic Regression(LR) and Outliers Addressed Data-Set, After ONEHOT
TEST #87: Average Accuracy is 85.79%, Average TPR is 85.85%, Average FNR is 14.15%



Results of ML Models: (sorted by False Negative Rate(FNR)) (top 10)

Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
TEST #108: Average Accuracy is 87.57%, Average TPR is 87.62%, Average FNR is 12.38%

Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
TEST #116: Average Accuracy is 86.87%, Average TPR is 86.93%, Average FNR is 13.07%

Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
TEST #110: Average Accuracy is 86.68%, Average TPR is 86.74%, Average FNR is 13.26%

Using ML Model: Random Forest(RF) and Outliers Addressed Data-Set, After ONEHOT, Afte
TEST #111: Average Accuracy is 86.48%, Average TPR is 86.53%, Average FNR is 13.47%

Using ML Model: XG Boost(XGB) and Original Data-Set, After ONEHOT, After SMOTE:
TEST #118: Average Accuracy is 86.48%, Average TPR is 86.53%, Average FNR is 13.47%

Using ML Model: Random Forest(RF) and Normalized Data-Set, After SMOTE:
TEST #109: Average Accuracy is 86.09%, Average TPR is 86.14%, Average FNR is 13.86%

Using ML Model: XG Boost(XGB) and Normalized Data-Set, After SMOTE:
TEST #117: Average Accuracy is 85.98%, Average TPR is 86.03%, Average FNR is 13.97%
```

## Hypertuning:

```
1  #For Hyper-tuning and secondary validation, only three dataset variations will be used.
2  #"Original Data-Set, After ONEHOT, After SMOTE"
3  #"Outliers Addressed Data-Set, After ONEHOT, After SMOTE"
4  #"Original Data-Set, After SMOTE"
5  #
6  #Copying these over to easier variables for these two stages.
7  df1 = df_smote.copy()
8  df1.name = df_smote.name #For some reason df.copy() doesn't copy the "name" attribute.
9
10 df2 = df_onehot_smote.copy()
11 df2.name = df_onehot_smote.name #For some reason df.copy() doesn't copy the "name" attribu
12
13 df3 = df_onehot_outliers_addressed_smote.copy()
14 df3.name = df_onehot_outliers_addressed_smote.name #For some reason df.copy() doesn't copy
15
16 df_list = []
17 df_list = [df1, df2, df3]
18
19 classCol = "Target"
```

```
1  #Based on the initial results, for hyper-tuning, focusing on models RF, XGB, and LR.
2
3  #Since hyper-tuning takes the most time to run, of everything in this project,
4  #wrapping this section up in a "if" statement so I can turn it off for future top down run
5  #Once hypertuned values are retrieved, the output will be copied into a text block for ref
6  #perform_tuning = True
7  perform_tuning = False
8
9  #Models used for Hyper-tuning.
10 rf = RandomForestClassifier()
11 xgb = XGBClassifier()
12 lr = LogisticRegression(max_iter=10000)
13
14 #Parameter Options tried for Hyper-tuning.
15 params_rf = {'n_estimators':[100,200,300,400,500], 'min_samples_leaf':[5, 10, 15, 20, 25,
16 params_xgb = {'n_estimators': [100,200,300,400,500,600,700,800,900,1000], 'learning_rate':
17 params_lr = {'solver':['newton-cg', 'lbfgs', 'sag', 'saga', 'liblinear'], 'penalty':['l2']
18
19 folds = 20
20
21 if perform_tuning:
22   params_rf = {'n_estimators':[100,200,300,400,500], 'min_samples_leaf':[5, 10, 15, 20, 25
23   params_xgb = {'n_estimators': [100,200,300,400,500,600,700,800,900,1000], 'learning_rate
24   params_lr = {'solver':['newton-cg', 'lbfgs', 'sag', 'saga', 'liblinear'], 'penalty':['l2
25
26   for idx,df_not_used in enumerate(df_list):
27     this_df = df_list[idx]
28
29     feature_names=this_df.columns[this_df.columns != classCol]
```

```
30      #80% training and 20% test
31      X_train, X_test, Y_train, Y_test = train_test_split(this_df.loc[:, feature_names], thi
32
33      grid_rf = GridSearchCV(rf, param_grid=params_rf, cv=folds)
34      grid_rf.fit(X_train, Y_train)
35      print(f"Using Data-Set: {this_df.name}:")
36      print("Hyper-Tuned Parameters for Random Forest:", grid_rf.best_params_)
37      print("")
38
39      rs_xgb =  RandomizedSearchCV(xgb, param_distributions=params_xgb, cv=folds)
40      rs_xgb.fit(X_train, Y_train)
41      print(f"Using Data-Set: {this_df.name}:")
42      print("Hyper-Tuned Parameters for XGBoost:", rs_xgb.best_params_)
43      print("")
44
45      grid_lr = GridSearchCV(lr, param_grid=params_lr, cv=folds)
46      grid_lr.fit(X_train, Y_train)
47      print(f"Using Data-Set: {this_df.name}:")
48      print("Hyper-Tuned Parameters for Logistic Regression:", grid_lr.best_params_)
49      print("")
50
```

**Results from hyper-tuning:**

Using Data-Set: Original Data-Set, After SMOTE: Hyper-Tuned Parameters for Random Forest: {'min_samples_leaf': 10, 'n_estimators': 500}

Using Data-Set: Original Data-Set, After SMOTE: Hyper-Tuned Parameters for XGBoost: {'n_estimators': 1000, 'learning_rate': 0.3}

Using Data-Set: Original Data-Set, After SMOTE: Hyper-Tuned Parameters for Logistic Regression: {'C': 100, 'penalty': 'l2', 'solver': 'newton-cg'}

Using Data-Set: Original Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters for Random Forest: {'min_samples_leaf': 5, 'n_estimators': 400}

Using Data-Set: Original Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters for XGBoost: {'n_estimators': 500, 'learning_rate': 0.2}

Using Data-Set: Original Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters for Logistic Regression: {'C': 100, 'penalty': 'l2', 'solver': 'newton-cg'}

Using Data-Set: Outliers Addressed Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters for Random Forest: {'min_samples_leaf': 5, 'n_estimators': 200}

Using Data-Set: Outliers Addressed Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters for XGBoost: {'n_estimators': 400, 'learning_rate': 0.5}

Using Data-Set: Outliers Addressed Data-Set, After ONEHOT, After SMOTE: Hyper-Tuned Parameters
for Logistic Regression: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}

```
1 #Apply Hyper-tuned parameters.
2 #I know there is an automated way of applying the Hyper-tuned parameters from the output o
3 #grid search, but this works too.
4
5 rf_tuned_df1 = RandomForestClassifier(min_samples_leaf = 10, n_estimators = 500)
6 xgb_tuned_df1 = XGBClassifier(n_estimators = 1000, learning_rate = 0.3)
7 lr_tuned_df1 = LogisticRegression(max_iter=10000, C=100, penalty="l2", solver="newton-cg")
8
9 rf_tuned_df2 = RandomForestClassifier(min_samples_leaf = 5, n_estimators = 400)
10 xgb_tuned_df2 = XGBClassifier(n_estimators = 500, learning_rate = 0.2)
11 lr_tuned_df2 = LogisticRegression(max_iter=10000, C=100, penalty="l2", solver="newton-cg")
12
13 rf_tuned_df3 = RandomForestClassifier(min_samples_leaf = 5, n_estimators = 200)
14 xgb_tuned_df3 = XGBClassifier(n_estimators = 400, learning_rate = 0.5)
15 lr_tuned_df3 = LogisticRegression(max_iter=10000, C=10, penalty="l2", solver="liblinear")
16
```

## Secondary Validation: (After Hypertuning):

```
1 #This will start testing all over again.
2 results = {}
3 test_num = 0
4 folds = 20
5
6 for test in ["LR_DF1","LR_DF2","LR_DF3","LR_TUNED_DF1","LR_TUNED_DF2","LR_TUNED_DF3",
7             "RF_DF1","RF_DF2","RF_DF3","RF_TUNED_DF1","RF_TUNED_DF2","RF_TUNED_DF3",
8             "XGB_DF1","XGB_DF2","XGB_DF3","XGB_TUNED_DF1","XGB_TUNED_DF2","XGB_TUNED_DF3"
9   if test == "LR_DF1":
10     model = lr
11     test_name = "Logistic Regression(LR)"
12     this_df = df1
13   elif test == "LR_DF2":
14     model = lr
15     test_name = "Logistic Regression(LR)"
16     this_df = df2
17   elif test == "LR_DF3":
18     model = lr
19     test_name = "Logistic Regression(LR)"
20     this_df = df3
21
22   elif test == "LR_TUNED_DF1":
23     model = lr_tuned_df1
24     test_name = "Logistic Regression(LR) - Tuned"
25     this_df = df1
26   elif test == "LR_TUNED_DF2":
27     model = lr_tuned_df2
```

```
28       test_name = "Logistic Regression(LR) - Tuned"
29       this_df = df2
30   elif test == "LR_TUNED_DF3":
31       model = lr_tuned_df3
32       test_name = "Logistic Regression(LR) - Tuned"
33       this_df = df3
34
35   elif test == "RF_DF1":
36       model = rf
37       test_name = "Random Forest(RF)"
38       this_df = df1
39   elif test == "RF_DF2":
40       model = rf
41       test_name = "Random Forest(RF)"
42       this_df = df2
43   elif test == "RF_DF3":
44       model = rf
45       test_name = "Random Forest(RF)"
46       this_df = df3
47
48   elif test == "RF_TUNED_DF1":
49       model = rf_tuned_df1
50       test_name = "Random Forest(RF) - Tuned"
51       this_df = df1
52   elif test == "RF_TUNED_DF2":
53       model = rf_tuned_df2
54       test_name = "Random Forest(RF) - Tuned"
55       this_df = df2
56   elif test == "RF_TUNED_DF3":
57       model = rf_tuned_df3
58       test_name = "Random Forest(RF) - Tuned"
59       this_df = df3
60
61   elif test == "XGB_DF1":
62       model = xgb
63       test_name = "XG Boost(XGB)"
64       this_df = df1
65   elif test == "XGB_DF2":
66       model = xgb
67       test_name = "XG Boost(XGB)"
68       this_df = df2
69   elif test == "XGB_DF3":
70       model = xgb
71       test_name = "XG Boost(XGB)"
72       this_df = df3
73
74   elif test == "XGB_TUNED_DF1":
75       model = xgb_tuned_df1
76       test_name = "XG Boost(XGB) - Tuned"
77       this_df = df1
78   elif test == "XGB_TUNED_DF2":
```

```
79     model = xgb_tuned_df2
80     test_name = "XG Boost(XGB) - Tuned"
81     this_df = df2
82   elif test == "XGB_TUNED_DF3":
83     model = xgb_tuned_df3
84     test_name = "XG Boost(XGB) - Tuned"
85     this_df = df3
86
87   doCV(test_num,test_name,model,this_df,folds,classCol,1)
88   test_num += 1
```

```
    For TEST #0: Logistic Regression(LR)              After Cross-Val, Using Data-Set: Origin
    For TEST #1: Logistic Regression(LR)              After Cross-Val, Using Data-Set: Origin
    For TEST #2: Logistic Regression(LR)              After Cross-Val, Using Data-Set: Outlie
    For TEST #3: Logistic Regression(LR) - Tuned      After Cross-Val, Using Data-Set: Origin
    For TEST #4: Logistic Regression(LR) - Tuned      After Cross-Val, Using Data-Set: Origin
    For TEST #5: Logistic Regression(LR) - Tuned      After Cross-Val, Using Data-Set: Outlie
    For TEST #6: Random Forest(RF)                    After Cross-Val, Using Data-Set: Origin
    For TEST #7: Random Forest(RF)                    After Cross-Val, Using Data-Set: Origin
    For TEST #8: Random Forest(RF)                    After Cross-Val, Using Data-Set: Outlie
    For TEST #9: Random Forest(RF) - Tuned            After Cross-Val, Using Data-Set: Origin
    For TEST #10: Random Forest(RF) - Tuned           After Cross-Val, Using Data-Set: Origi
    For TEST #11: Random Forest(RF) - Tuned           After Cross-Val, Using Data-Set: Outli
    For TEST #12: XG Boost(XGB)                       After Cross-Val, Using Data-Set: Origi
    For TEST #13: XG Boost(XGB)                       After Cross-Val, Using Data-Set: Origi
    For TEST #14: XG Boost(XGB)                       After Cross-Val, Using Data-Set: Outli
    For TEST #15: XG Boost(XGB) - Tuned               After Cross-Val, Using Data-Set: Origi
    For TEST #16: XG Boost(XGB) - Tuned               After Cross-Val, Using Data-Set: Origi
    For TEST #17: XG Boost(XGB) - Tuned               After Cross-Val, Using Data-Set: Outli
```

```
 1 #View these results, after hyper-tuning, just the top 5 using each sort method:
 2 range_limit = min(5,results[f"LAST_TEST"]+1) #Top 5 results desired.
 3
 4 results_list = list(range(0,results[f"LAST_TEST"]+1))
 5 results_list.sort(key=getTestAccuracy, reverse=True)
 6
 7 print("Results of ML Models: (sorted by accuracy) (top 5)")
 8 print("")
 9 for i in results_list[0:range_limit]:  #Count starts at zero, so dont need to add one.
10   displayResult(i)
11 print("")
12 print("")
13
14
15 results_list = list(range(0,results[f"LAST_TEST"]+1))
16 results_list.sort(key=getTestFNR, reverse=False)
17
18 print("Results of ML Models: (sorted by False Negative Rate(FNR)) (top 5)")
19 print("")
20 for i in results_list[0:range_limit]:  #Count starts at zero, so dont need to add one.
21   displayResult(i)
```

```
22 print("")
23 print("")
24
25
26 #View these results, after hyper-tuning, just the bottom 5 using each sort method:
27 range_limit = min(5,results[f"LAST_TEST"]+1) #Bottom 5 results desired.
28
29 results_list = list(range(0,results[f"LAST_TEST"]+1))
30 results_list.sort(key=getTestAccuracy, reverse=False)
31
32 print("Results of ML Models: (sorted by accuracy) (bottom 5)")
33 print("")
34 for i in results_list[0:range_limit]:  #Count starts at zero, so dont need to add one.
35   displayResult(i)
36 print("")
37 print("")
38
39
```

        Results of ML Models: (sorted by accuracy) (top 5)

        Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
        TEST #6: Average Accuracy is 88.44%, Average TPR is 88.37%, Average FNR is 11.63%

        Using ML Model: Random Forest(RF) - Tuned and Original Data-Set, After ONEHOT, After
        TEST #10: Average Accuracy is 88.24%, Average TPR is 88.18%, Average FNR is 11.82%

        Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
        TEST #12: Average Accuracy is 88.24%, Average TPR is 88.17%, Average FNR is 11.83%

        Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
        TEST #7: Average Accuracy is 88.14%, Average TPR is 88.09%, Average FNR is 11.91%

        Using ML Model: Random Forest(RF) and Outliers Addressed Data-Set, After ONEHOT, Afte
        TEST #8: Average Accuracy is 87.85%, Average TPR is 87.78%, Average FNR is 12.22%


        Results of ML Models: (sorted by False Negative Rate(FNR)) (top 5)

        Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
        TEST #6: Average Accuracy is 88.44%, Average TPR is 88.37%, Average FNR is 11.63%

        Using ML Model: Random Forest(RF) - Tuned and Original Data-Set, After ONEHOT, After
        TEST #10: Average Accuracy is 88.24%, Average TPR is 88.18%, Average FNR is 11.82%

        Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
        TEST #12: Average Accuracy is 88.24%, Average TPR is 88.17%, Average FNR is 11.83%

        Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
        TEST #7: Average Accuracy is 88.14%, Average TPR is 88.09%, Average FNR is 11.91%

        Using ML Model: Random Forest(RF) and Outliers Addressed Data-Set, After ONEHOT, Afte
        TEST #8: Average Accuracy is 87.85%, Average TPR is 87.78%, Average FNR is 12.22%

Results of ML Models: (sorted by accuracy) (bottom 5)

Using ML Model: Logistic Regression(LR) - Tuned and Original Data-Set, After SMOTE:
TEST #3: Average Accuracy is 85.30%, Average TPR is 85.19%, Average FNR is 14.81%

Using ML Model: Logistic Regression(LR) and Original Data-Set, After SMOTE:
TEST #0: Average Accuracy is 85.30%, Average TPR is 85.19%, Average FNR is 14.81%

Using ML Model: XG Boost(XGB) - Tuned and Outliers Addressed Data-Set, After ONEHOT,
TEST #17: Average Accuracy is 86.17%, Average TPR is 86.10%, Average FNR is 13.90%

Using ML Model: XG Boost(XGB) - Tuned and Original Data-Set, After SMOTE:
TEST #15: Average Accuracy is 86.36%, Average TPR is 86.28%, Average FNR is 13.72%

Using ML Model: Logistic Regression(LR) - Tuned and Original Data-Set, After ONEHOT,
TEST #4: Average Accuracy is 86.37%, Average TPR is 86.28%, Average FNR is 13.72%

```
 1 #Comparing Results.
 2 #Due to random splits, the results change from run to run, but for the most part XGBoost w
 3 #whereas Random Forest produces very similar results tuned, or not tuned.  Consistently th
 4 #
 5 #
 6 #Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
 7 #TEST #6: Average Accuracy is 88.44%, Average TPR is 88.37%, Average FNR is 11.63%
 8
 9 #Using ML Model: Random Forest(RF) - Tuned and Original Data-Set, After ONEHOT, After SMOT
10 #TEST #10: Average Accuracy is 88.24%, Average TPR is 88.18%, Average FNR is 11.82%
11
12 #Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
13 #TEST #12: Average Accuracy is 88.24%, Average TPR is 88.17%, Average FNR is 11.83%
14
15 #Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
16 #TEST #7: Average Accuracy is 88.14%, Average TPR is 88.09%, Average FNR is 11.91%
17
18 #
19 # But are the differences between models, and data-sets statistically significant?
20 #
21 # For example, we have these three pairs of results:
22 #
23 # Pair #1: Same Data-set/Different models.
24 #Using ML Model: Random Forest(RF) and Original Data-Set, After SMOTE:
25 #TEST #6: Average Accuracy is 88.44%, Average TPR is 88.37%, Average FNR is 11.63%
26 #
27 #Using ML Model: XG Boost(XGB) and Original Data-Set, After SMOTE:
28 #TEST #12: Average Accuracy is 88.24%, Average TPR is 88.17%, Average FNR is 11.83%
29 #
```

```
30 #
31 # Pair #2: Same Data-set/RF vs RF Tuned
32 #Using ML Model: Random Forest(RF) - Tuned and Original Data-Set, After ONEHOT, After SMOT
33 #TEST #10: Average Accuracy is 88.24%, Average TPR is 88.18%, Average FNR is 11.82%
34
35 #Using ML Model: Random Forest(RF) and Original Data-Set, After ONEHOT, After SMOTE:
36 #TEST #7: Average Accuracy is 88.14%, Average TPR is 88.09%, Average FNR is 11.91%
37 #
38 # Each pair uses the same data-set, which make them suitable to compare results against
39 # each other using a paired dependant t-test.
40 #
41 #
42 def my_ttest(x,y,verbose=0):
43   t,p = stats.ttest_rel(getTestAccVec(x),getTestAccVec(y))
44   x_name = results[f"TEST{x}#NAME"]
45   x_df_name = results[f"TEST{x}#DFNAME"]
46   y_name = results[f"TEST{y}#NAME"]
47   y_df_name = results[f"TEST{y}#DFNAME"]
48
49   if p >= 0.05:
50     if verbose > 0:
51       print(f"TEST #{x}: {x_name}: Using Data-Set: {x_df_name} &")
52       print(f"TEST #{y}: {y_name}: Using Data-Set: {y_df_name}")
53       print(f"Results for Accuracy are NOT significantly different from each other, with t
54       print("")
55     return(False)
56   else:
57     if verbose > 0:
58       print(f"TEST #{x}: {x_name}: Using Data-Set: {x_df_name} &")
59       print(f"TEST #{y}: {y_name}: Using Data-Set: {y_df_name}")
60       print(f"Results for Accuracy are significantly different from each other, with t={t}
61       print("")
62     return(True)
63
64 #Pair #1:
65 res = my_ttest(6,12,1)
66
67 #Pair #2:
68 res = my_ttest(10,7,1)
```

```
    TEST #6: Random Forest(RF): Using Data-Set: Original Data-Set, After SMOTE &
    TEST #12: XG Boost(XGB): Using Data-Set: Original Data-Set, After SMOTE
    Results for Accuracy are NOT significantly different from each other, with t=0.26184748

    TEST #10: Random Forest(RF) - Tuned: Using Data-Set: Original Data-Set, After ONEHOT, Af
    TEST #7: Random Forest(RF): Using Data-Set: Original Data-Set, After ONEHOT, After SMOTE
    Results for Accuracy are NOT significantly different from each other, with t=0.19792316
```

```
 1 #The paired dependant t-test above revealed there is no significant difference between any
```