

Problem Statement

The Case Study "Predictive Equipment Failures" deals with prediction of the failure of an equipment before it can fail. As an employee while working in an MNC, we had been given the task of building IoT solutions for the underlying infrastructure of the company. The entire IoT system was built with sensors and connected over the Internet to Cloud to monitor readings, control devices such as HVAC, Diesel Generators etc. The import diesel generator fail, the cost of repair at that instant in an emergency is multiple times of the cost that it would take to repair when done as maintenance. Such costs bring down revenue for the company even though the company had the data before hand but was not able to convert the same into useful information.

This case study exactly deals with a similar scenario as posed by a company named "ConocoPhillips" in West Texas. The company has several oil wells classified as stripper wells. Stripper wells produce low volumes of oil at well level, but at a country level of the US, they amount for significant domestic oil production. Since the wells have low throughput, if the equipment for extracting the oil fails, then profit margin too is impacted se

ML Formulation

The Task is to predict the failure of an equipment in advance which makes this problem a Binary Classification Problem where "1" represents down hole equipment failure and "0" represent no failure.

Business Constraints

- We should not miss any failures since it would lead to high costs.
- There is no time constraint, but the model should take a few seconds to a minute for prediction but not hours.
- Interpretability is not important. But if available it would be useful in identifying the exact part where the equipment might fail in future.

Dataset Analysis

The dataset provided "ConocoPhillips" consists of data from 107 sensors at surface level and down-hole equipment. The dataset consists of two types of sensor columns namely, a. Measure columns – single measurement of the sensor b. Histogram bin columns – a set of 10 columns with different bins of a sensor that show its distribution over time.

Thus, in total there are 100 measure column sensors and 7 histogram bin column sensor values resulting in 170 columns of only sensors. The target value consists of "1" for equipment failure and "0" for no failure with majority of the target being "0" and very few values being "1" thus resulting in a highly imbalanced dataset. There are a total of 60,000 rows in the dataset.

Performance Metric

Since the dataset is highly imbalanced, the ideal performance metric would be precision, recall and F1 scores especially not to miss failure cases. Other performance metrics like accuracy and ROC-AUC would provide a significantly high value due to the imbalance and not provide a clear picture for the minority class.

Among F1 score, we will choose the macro averaged f1 score since we would like to treat both class as equals. Micro Average will be used when we want to maximize the classification of a particular class which is not the case in this problem study. Due to imbalance dataset, we will get a high micro average f1 score which is not called for.

<https://datascience.stackexchange.com/questions/36862/macro-or-micro-average-for-imbalanced-class-problems>

Load Data

```
In [87]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import random

from sklearn.manifold import TSNE
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.impute import KNNImputer
from sklearn.decomposition import TruncatedSVD
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import plot_confusion_matrix
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import StackingClassifier
from sklearn.ensemble import VotingClassifier
from sklearn.tree import DecisionTreeRegressor
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from prettytable import PrettyTable
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from joblib import dump, load
import pickle

In [104]: data_df = pd.read_csv('equip_failures_training_set.csv')
data_df.head()
data_df_copy = data_df.copy()
```

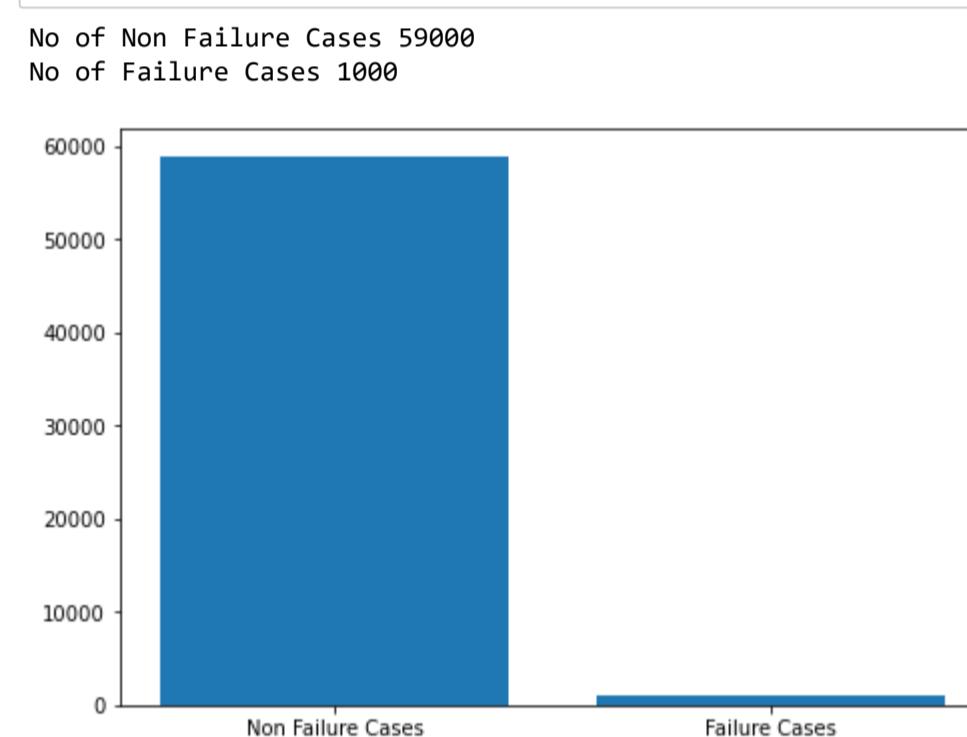
Exploratory Data Analysis

Count the imbalanced dataset target values

```
In [29]: print(data_df.shape)
(60000, 172)

In [30]: # Obtaining the total count of each class in the target variable
targets = data_df['target'].value_counts()

figure = plt.figure()
axes = figure.add_axes([0,0,1,1])
target_values = ['Non Failure Cases', 'Failure Cases']
counts = [targets[0], targets[1]]
print("No of Non Failure Cases",targets[0])
print("No of Failure Cases",targets[1])
axes.bar(target_values,counts)
plt.show()
```



We see that the data is highly imbalanced with only 1000 data points in the minority class.

PDF plots of each non temporal sensor

```
In [31]: sensor_column_names = []
for column in data_df.columns:
    if column != 'id' and column != 'target':
        column_name = column.split('_')
        if column_name[1] == 'histogram':
            sensor_column_names.append(column)

In [32]: for col in data_df.columns:
    data_df[col] = pd.to_numeric(data_df[col], errors='coerce')

In [33]: df_0 = data_df[data_df['target']==0]
df_1 = data_df[data_df['target']==1]

In [16]: fig, axes = plt.subplots(25, 4, figsize=(20,80))
axes = axes.flatten()

for i,column in enumerate(sensor_column_names):
    sns.distplot(df_0[column], hist=False, ax=axes[i], label='0')
    sns.distplot(df_1[column], hist=False, ax=axes[i], label='1')
    axes[i].set_title(column)
    fig.tight_layout()

c:\users\manish\dalvi\appdata\local\programs\python\python37\lib\site-packages\seaborn\distributions.py:283: UserWarning: Data must have variance to compute a kernel density estimate.
warnings.warn(msg, UserWarning)
c:\users\manish\dalvi\appdata\local\programs\python\python37\lib\site-packages\seaborn\distributions.py:283: UserWarning: Data must have variance to compute a kernel density estimate.
warnings.warn(msg, UserWarning)
```

Comments

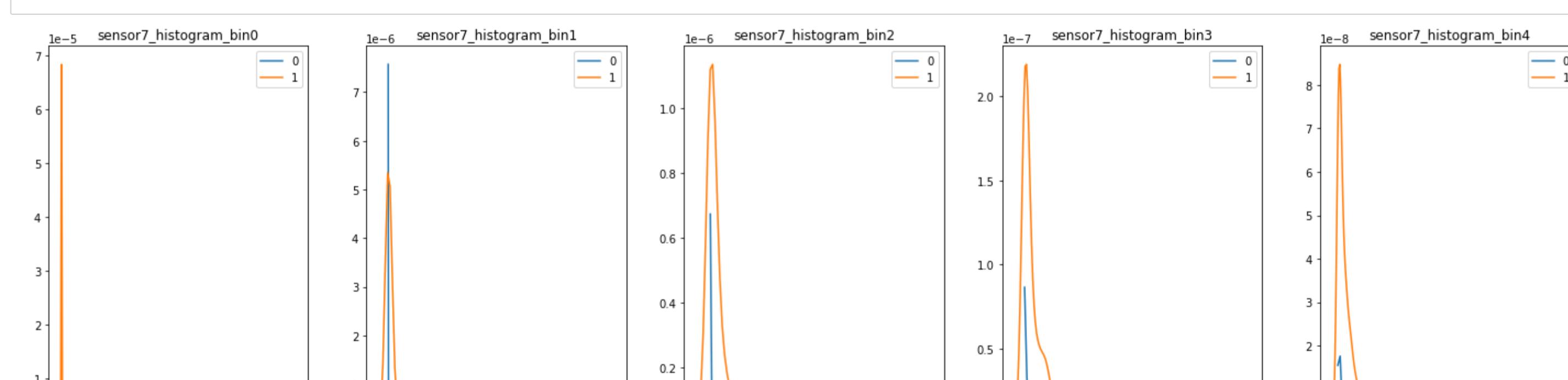
- Most of the graphs are overlapping due to which it is difficult to segregate the data points by a particular sensor.
- There are few sensors whose graph is visually separable such as Sensor 1, Sensor 8, Sensor 14-17, Sensor 27, Sensor 67 and several more.
- These separable distributions add greater value to the final classification hence it would be necessary to look at their importance once we apply the same on Models.

Histogram plots

```
In [34]: histogram_column_names = []
for column in data_df.columns:
    if column != 'id' and column != 'target':
        column_name = column.split('_')
        if column_name[1] == 'histogram':
            histogram_column_names.append(column)

In [388]: fig, axes = plt.subplots(14, 5, figsize=(20,80))
axes = axes.flatten()

for i,column in enumerate(histogram_column_names):
    sns.distplot(df_0[column], hist=False, ax=axes[i], label='0')
    sns.distplot(df_1[column], hist=False, ax=axes[i], label='1')
    axes[i].set_title(column)
    fig.tight_layout()
```



Comments

- The histogram bins look similar to power law graphs.
- Most of the values are near zero suggesting the high number of 0 values in the bins.
- After analyzing all the graphs, it is difficult to visually separate the values depending on their distributions of bins over time.

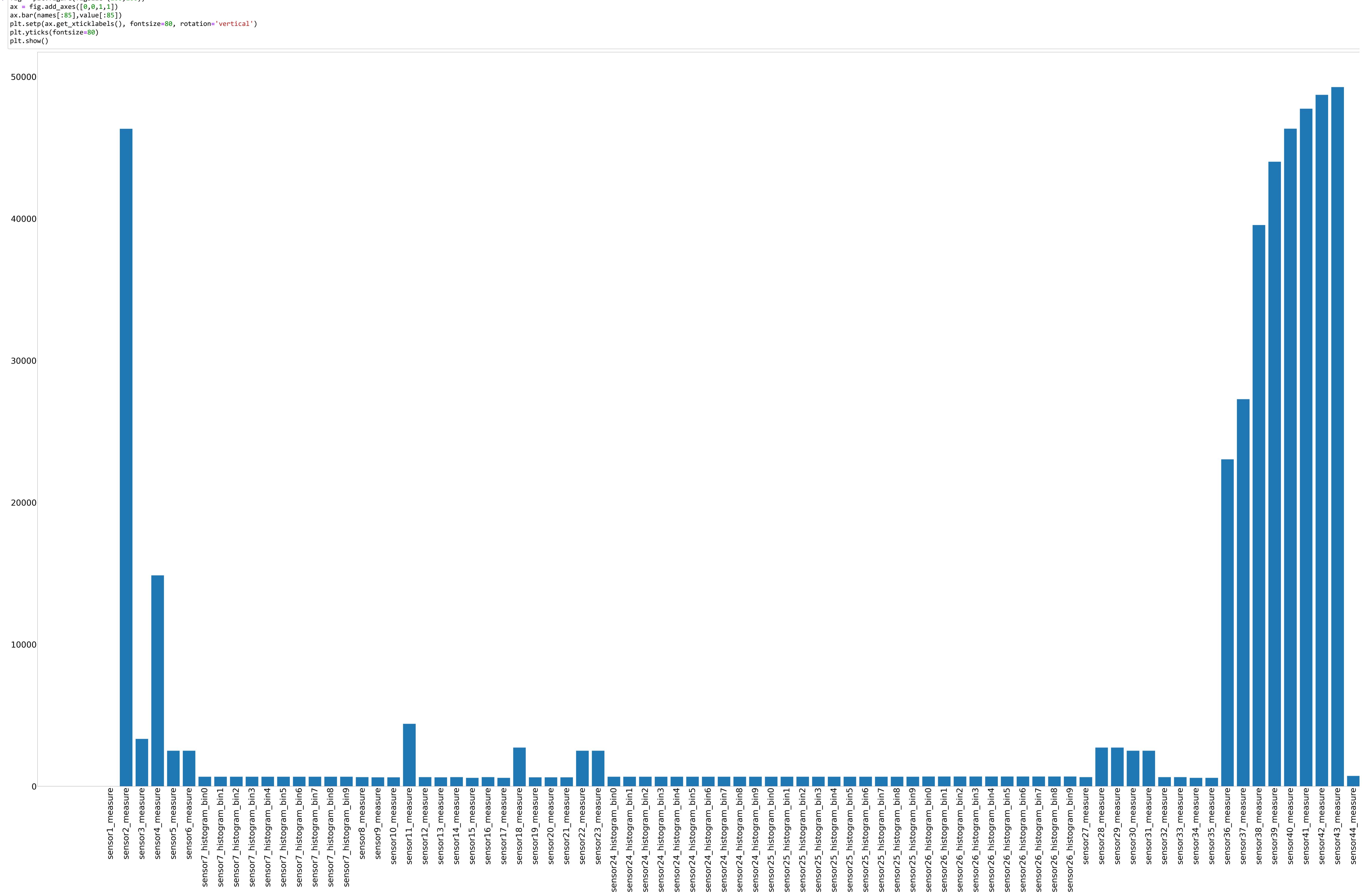
Check Null values in each feature

```
In [35]: # Checking null values in each column
null_values = data_df.isnull().sum(axis = 0)
```

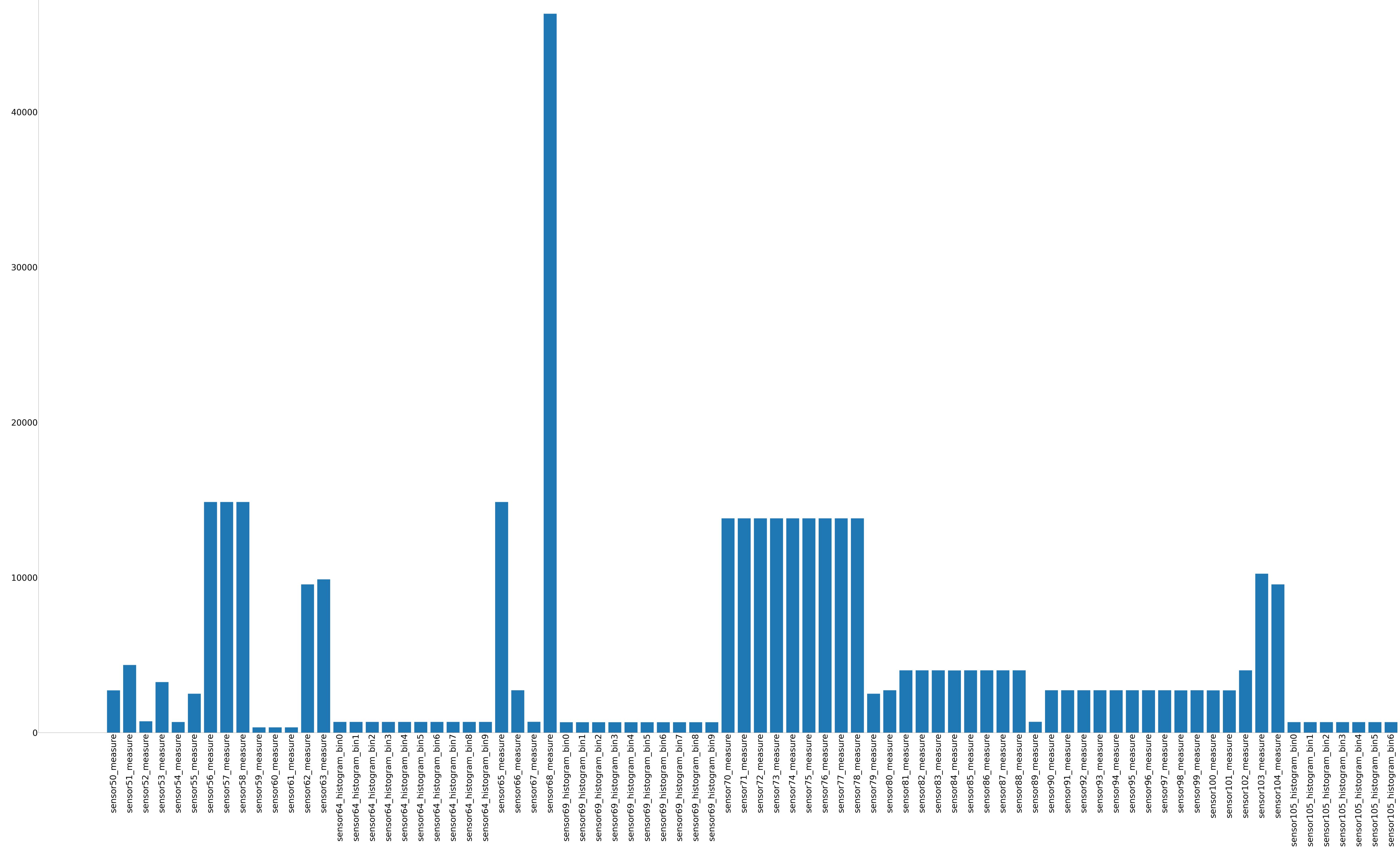
```
In [36]: value = []
names = []

for k, v in null_values.iteritems():
    if k != 'target' and k != 'id':
        names.append(k)
        value.append(v)

print(len(names), len(value))
```



```
In [313]: fig = plt.figure(figsize=(200,100))
ax = fig.add_axes([0,0,1,1])
ax.bar(names[85:170],value[85:170])
plt.setp(ax.get_xticklabels(), fontsize=80, rotation='vertical')
plt.yticks(fontsize=80)
plt.show()
```



Comments

- There are few sensors who have most of the values are NULL namely Sensor 2, Sensor 38-43 and sensor 68
 - These sensors columns will not contribute much to the final output and performing imputation will also spoil the nature of the data.

Dropping columns with more than 50% Null values in the rows

We can modify the data frame in this situation before train-test split it because we are removing the entire column and not performing any modification to value or implementation.

We can modify the dataframe in this situation before train, test split is because we are removing the entire column and not performing any modification to value or imputation.

In [37]:

```
null_values = null_values.sort_values(ascending=False)
```

Before Dropping Null values, checking their correlation with the target values for each column.

In [38]:

```
null_values_list = []
for column_name, value in null_values.items():
    if value > 30000:
        null_values_list.append(column_name)

print(null_values_list)
```

```
['sensor43 measure', 'sensor42 measure', 'sensor41 measure', 'sensor40 measure', 'sensor68 measure', 'sensor2 measure', 'sensor39 measure', 'sensor38 measure']
```

```
In [39]: for column in null_values_list:
    print("Correlation between " + column + " and target is: " + str(data_df[column].corr(data_df['target'])))

Correlation between sensor43_measure and target is: -0.14302546793845026
Correlation between sensor42_measure and target is: -0.1316889144598184
Correlation between sensor41_measure and target is: -0.11065586976877181
Correlation between sensor40_measure and target is: -0.086662692353530962
Correlation between sensor39_measure and target is: 0.031922769466329883
Correlation between sensor38_measure and target is: -0.057736863918924256
Correlation between sensor37_measure and target is: -0.021159379389264103
```

We see no correlation between the columns with null values >50% and the target column, thus we can drop the specified columns

```
In [40]: # If the number of null values are greater than 50% of the data set it is better to drop them since they would not contribute anything.
print("Number of Columns before dropping null columns",data_df.shape[1])
for column_name, value in null_values.items():
    if value > 30000:
        data_df = data_df.drop([column_name], axis=1)
print("Number of Columns after dropping null columns",data_df.shape[1])
```

Check for 0's in each feature

```
In [41]: # Dropping columns with 95% zeros
zero_columns = []
for column in data_df.columns:
    if column != 'target' and column != 'id':
        value = data_df[column].value_counts()
        if 0 in value:
            if value[0] >= (0.95*60000):
                zero_columns.append(column)
                print(column, '--- % of zeros', (value[0]*100)/60000)

sensor7_histogram_bin0 --- % of zeros 98.555
sensor7_histogram_bin1 --- % of zeros 97.645
sensor19_measure --- % of zeros 98.91666666666667
sensor21_measure --- % of zeros 98.85
sensor24_histogram_bin0 --- % of zeros 97.939
sensor24_histogram_bin1 --- % of zeros 96.98666666666666
sensor24_histogram_bin2 --- % of zeros 96.98666666666666
sensor24_histogram_bin3 --- % of zeros 96.98666666666666
sensor24_histogram_bin4 --- % of zeros 95.26333333333334
sensor24_histogram_bin5 --- % of zeros 97.98166666666667
sensor25_histogram_bin9 --- % of zeros 95.74166666666666
sensor64_histogram_bin0 --- % of zeros 95.5
sensor106_measure --- % of zeros 95.035
```

Checking the correlation between the columns with >95% zeros and the target column

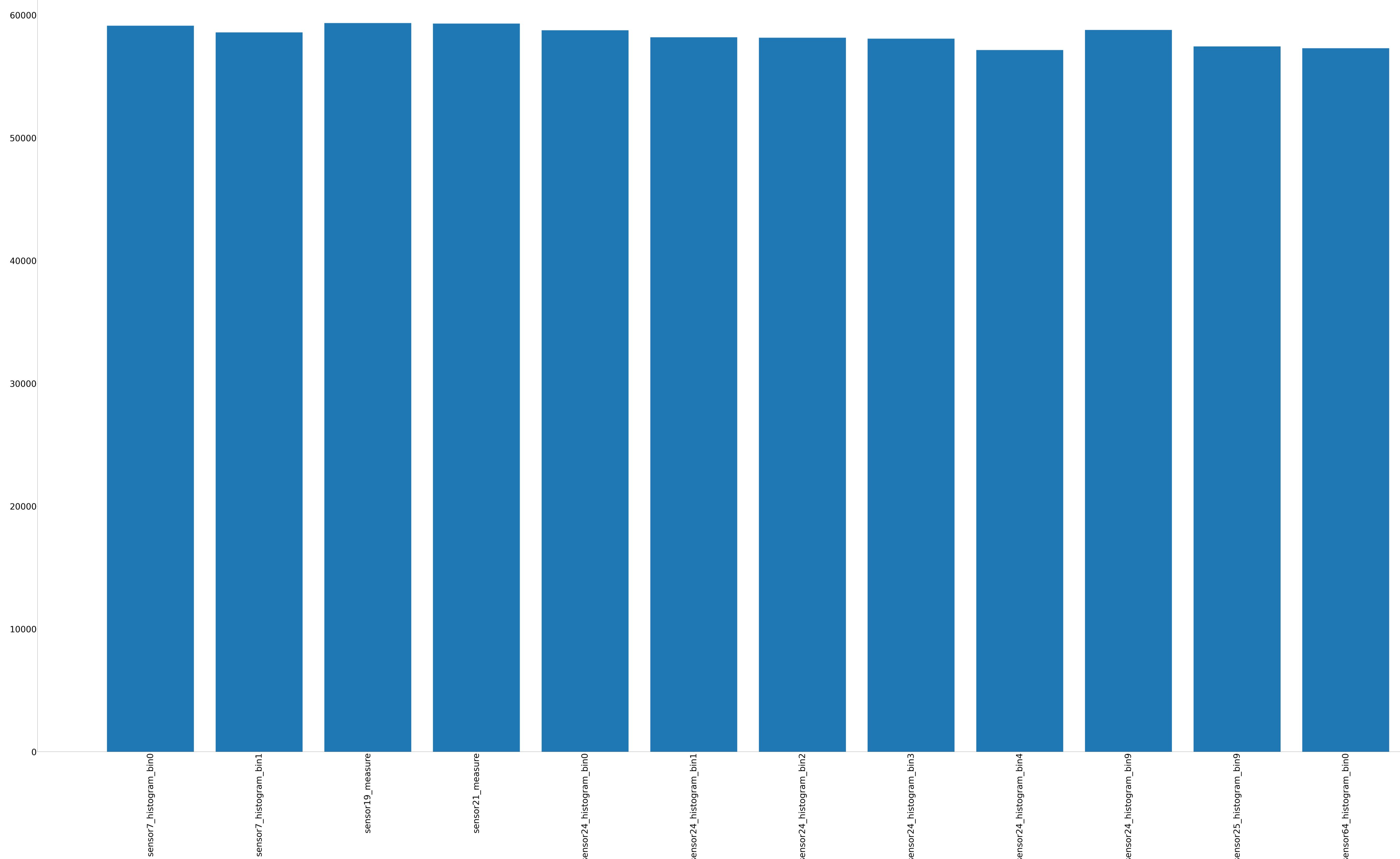
```
In [42]: for column in zero_columns:
    print("Correlation between " + column + " and target is: " + str(data_df[column].corr(data_df['target'])))

Correlation between sensor7_histogram_bin0 and target is: 0.01236794996627984
Correlation between sensor7_histogram_bin1 and target is: 0.1914529489338661
Correlation between sensor19_measure and target is: 0.044101519644678195
Correlation between sensor21_measure and target is: 0.09313679837210851
Correlation between sensor24_histogram_bin0 and target is: 0.1579831212668717
Correlation between sensor24_histogram_bin1 and target is: 0.16556112873662185
Correlation between sensor24_histogram_bin2 and target is: 0.16556112873662185
Correlation between sensor24_histogram_bin3 and target is: 0.19103827667736897
Correlation between sensor24_histogram_bin4 and target is: 0.1799881213531043
Correlation between sensor24_histogram_bin5 and target is: 0.042447918598523156
Correlation between sensor25_histogram_bin9 and target is: 0.1322741693510897
Correlation between sensor106_measure and target is: 0.018399967244495879
```

In [43]: **## We see no significant correlation between the columns with >95% zeros and the target columns, and hence we can drop the same.**

```
In [44]: value = []
for column in zero_columns:
    value.append(data_df[column].value_counts()[0])

In [318]: fig = plt.figure(figsize=(200,100))
ax = fig.add_axes([0,0,1,1])
ax.bar(zero_columns,value)
plt.setp(ax.get_xticklabels(), fontsize=80, rotation='vertical')
plt.yticks(fontsize=80)
plt.show()
```



Comments

- There are few sensors who have most of the values are Zero namely Sensor 7 bin 0-1, Sensor 10, Sensor 21, Sensor 24 bin 0-4, Sensor 24 bin9, Sensor 25 bin9, Sensor 64 bin0 and sensor 106 measure.
- These sensors have more than 95% of their values to zero thus providing no contribution to final classification.

Dropping columns with >95% number of rows of zero

We can modify the data frame in this situation before train, test split is because we are removing the entire column and not performing any modification to value or imputation.

```
In [45]: print("Number of columns before dropping Zero columns:", data_df.shape[1])
data_df = data_df.drop(zero_columns, axis=1)
print("Number of columns after dropping Zero columns:", data_df.shape[1])
```

Number of columns before dropping Zero columns: 164
Number of columns after dropping Zero columns: 151

Check for distinction between the two classes using quantile and box plots

```
In [46]: df_0 = data_df[data_df['target']==0]
df_1 = data_df[data_df['target']==1]

In [47]: # Checking the percentile of 0 and 1 of each column to see which features provide better separation.
# comparing if 90th percentile of 0 < 10th percentile of 1 and also 10th percentile of 1 > 90th percentile of 0
# Choosing Columns between 10 and 90th percentile with significant difference in the two class distribution

quantiles_0 = df_0.quantile([0.1, 0.9], axis=0)
quantiles_1 = df_1.quantile([0.1, 0.9], axis=0)

important_features = []
for col in quantiles_0.columns:
    if col != 'id' and col != 'target':
        col_name = col.split("-")
        if (quantiles_0[col].iloc[0] > quantiles_1[col].iloc[1]) or (quantiles_1[col].iloc[0] > quantiles_0[col].iloc[1]) and (col_name[1] != 'histogram'):
            important_features.append(col)

# print(len(important_features))
```

These are the features which may add importance in the decision making at the end since the 10th and 90th percentile of each of these features are separable and hence useful in decision making. The sensors that are part of this are as follows

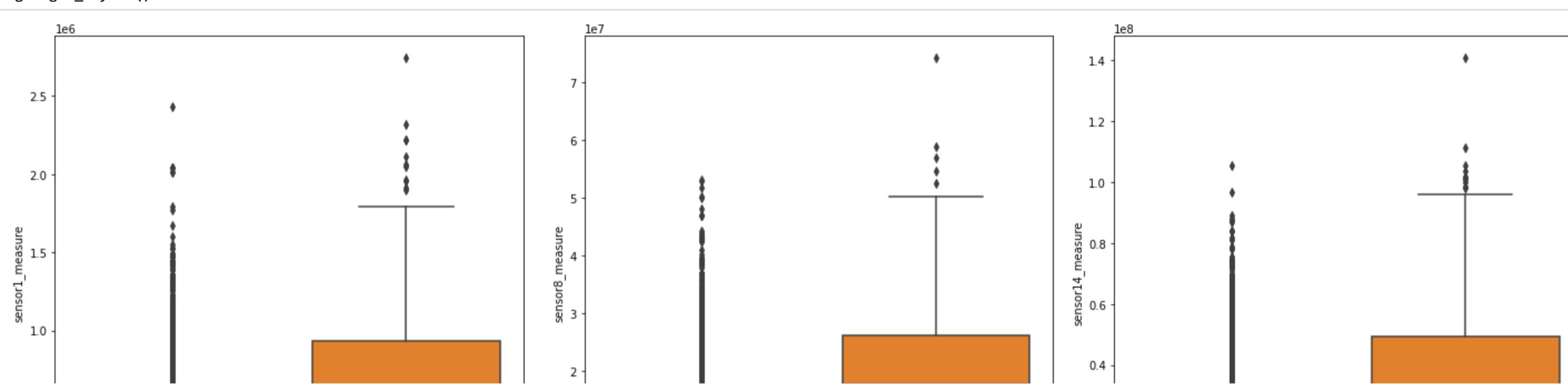
```

sensor1_measure
sensor2_measure
sensor3_measure
sensor4_measure
sensor5_measure
sensor6_measure
sensor7_measure
sensor8_measure
sensor9_measure
sensor10_measure
sensor11_measure
sensor12_measure
sensor13_measure
sensor14_measure
sensor15_measure
sensor16_measure
sensor17_measure
sensor18_measure
sensor19_measure
sensor20_measure
sensor21_measure
sensor22_measure
sensor23_measure
sensor24_measure
sensor25_measure
sensor26_measure
sensor27_measure
sensor28_measure
sensor29_measure
sensor30_measure
sensor31_measure
sensor32_measure
sensor33_measure
sensor34_measure
sensor35_measure
sensor36_measure
sensor37_measure
sensor38_measure
sensor39_measure
sensor40_measure
sensor41_measure
sensor42_measure
sensor43_measure
sensor44_measure
sensor45_measure
sensor46_measure
sensor47_measure
sensor48_measure
sensor49_measure
sensor50_measure
sensor51_measure
sensor52_measure
sensor53_measure
sensor54_measure
sensor55_measure
sensor56_measure
sensor57_measure
sensor58_measure
sensor59_measure
sensor60_measure
sensor61_measure
sensor62_measure
sensor63_measure
sensor64_measure
sensor65_measure
sensor66_measure
sensor67_measure
sensor68_measure
sensor69_measure

```

```
In [22]: fig, axes = plt.subplots(6, 3, figsize=(20,40))
axes = axes.flatten()

for i, column in enumerate(important_features):
    ax = sns.boxplot(x="target", y=column, data=data_df, ax=axes[i])
fig.tight_layout()
```



Comments

- The above box plots are for the features whose distribution is highly separable where the 10th percentile of one class is greater than 90th percentile of the other class.
- Since most of the data can be classified visually from these sensors, it is important to note that these features might contribute more to the final classification and also it would be important to note if these features are correlated.

TSNE plot for all the features

```
In [48]: for col in data_df_copy.columns:
    data_df_copy[col] = pd.to_numeric(data_df_copy[col], errors='coerce')

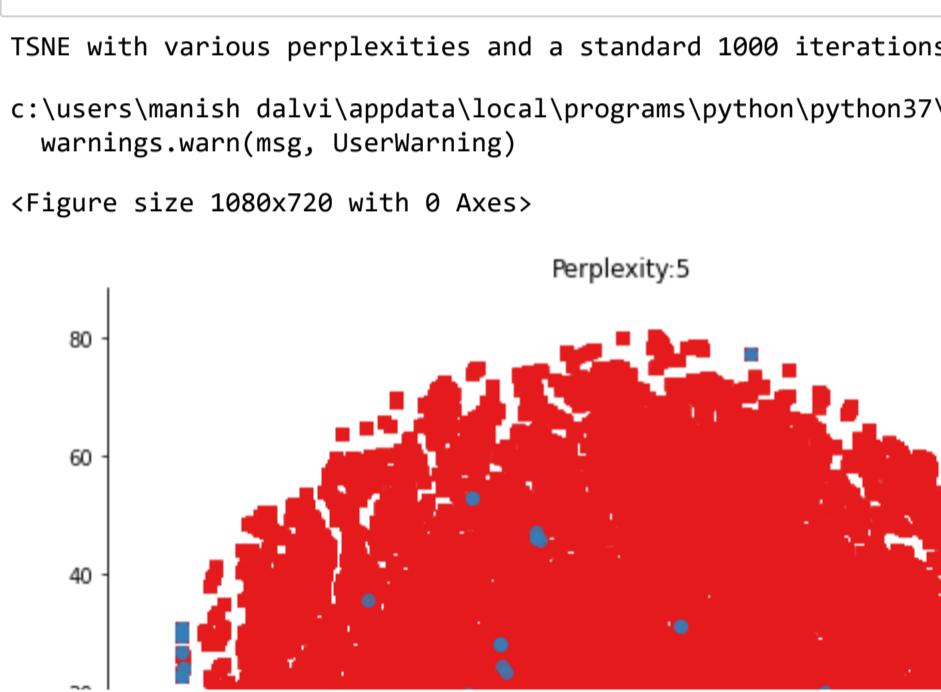
In [49]: # Using the copy of the dataset which was saved in the beginning since the other variable of dataset was modified.
data_df_copy = data_df_copy.fillna(data_df_copy.median())

y_true = data_df_copy['target']
data_df_copy = data_df_copy.drop(['target'],axis=1)
data_df_copy = data_df_copy.drop(['id'],axis=1)
```

```
In [376]: # Computing TSNE for each perplexity and displaying graph using seaborn plots
```

```
perplexity = [5, 10, 20, 40]
print("TSNE with various perplexities and a standard 1000 iterations")

for ppx in perplexity:
    tsne = TSNE(n_components=2, verbose=0, perplexity=ppx, n_iter=1000)
    tsne_data = tsne.fit_transform(data_df_copy)
    plt.figure(figsize=(15,10))
    df = pd.DataFrame({'x':tsne_data[:,0], 'y':tsne_data[:,1], 'label':y_true})
    sns.lmplot(data=df, x="x", y="y", hue="label", fit_reg=False, size=8, palette="Set1", markers=['s','o'])
    plt.title("Perplexity:{}".format(ppx))
    plt.show()
```



Comments

- From each of the TSNE plots we see that the data gets better grouped but the data is still overlapping on a 2D surface and thus converting the higher dimension data to a 2D data does not help in separation.

TSNE plot for important non time based sensors and easily separable features

```
In [50]: important_df = data_df_copy[['sensor1_measure','sensor8_measure', 'sensor14_measure', 'sensor15_measure', 'sensor16_measure', 'sensor17_measure', 'sensor27_measure', 'sensor32_measure', 'sensor33_measure', 'sensor34_measure', 'sensor35_measure', 'sensor46_measure', 'sensor47_measure', 'sensor59_measure', 'sensor61_measure', 'sensor62_measure']]
```

```
In [378]: # Computing TSNE for each perplexity and displaying graph using seaborn plots
```

```
perplexity = [5, 10, 20, 40]
for ppx in perplexity:
    tsne = TSNE(n_components=2, verbose=0, perplexity=ppx, n_iter=1000)
    tsne_data = tsne.fit_transform(important_df)
    plt.figure(figsize=(15,10))
    df = pd.DataFrame({'x':tsne_data[:,0], 'y':tsne_data[:,1], 'label':y_true})
    sns.lmplot(data=df, x="x", y="y", hue="label", fit_reg=False, size=8, palette="Set1", markers=['s','o'])
    plt.title("Perplexity:{}".format(ppx))
    plt.show()
```

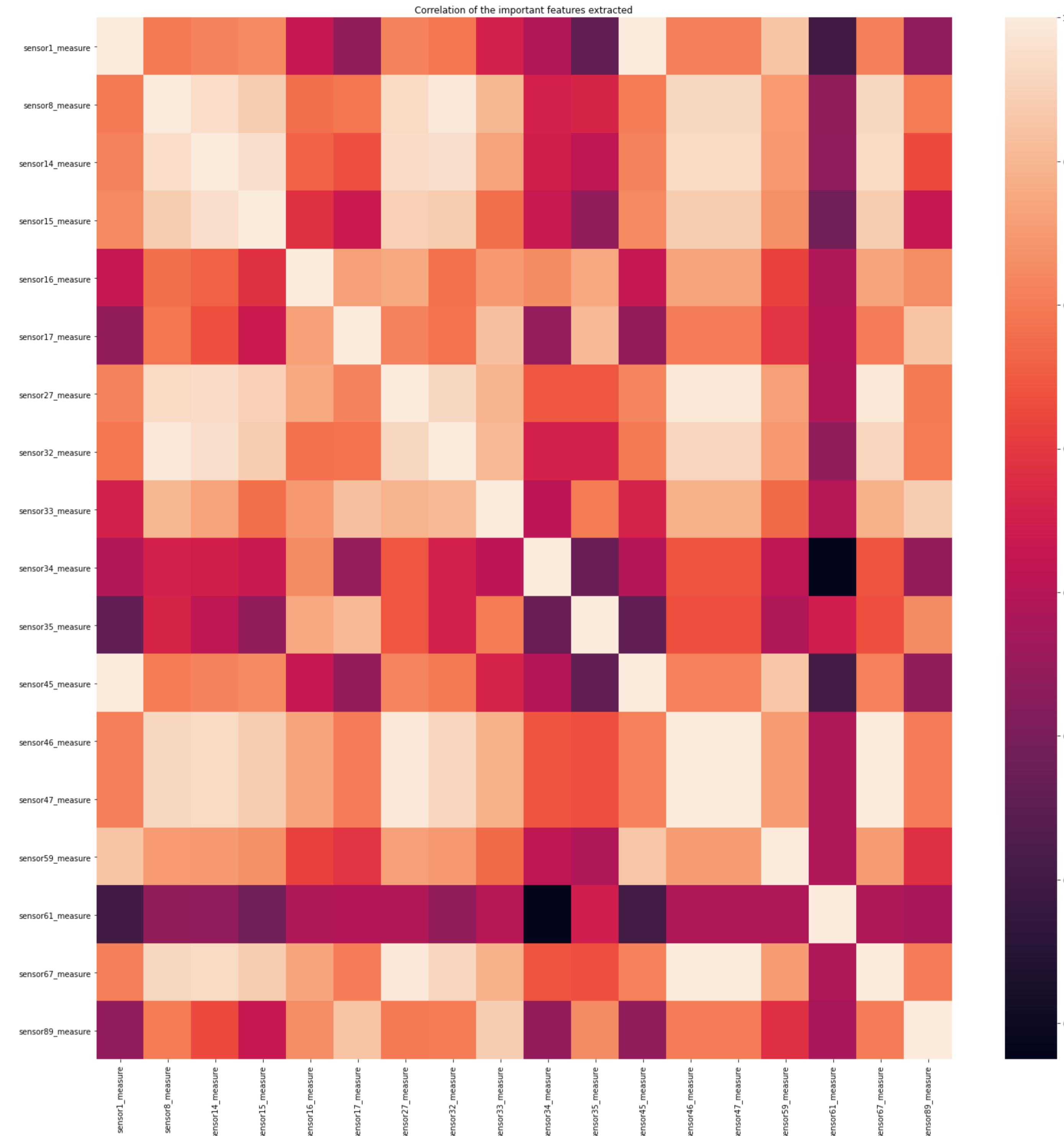


Comments

- Even with features that have their box plots mostly separable, we still see the overlap between the two classes and hence no defining boundary in 2D for their classification.

Correlation between easily separable features

```
In [51]: # Calculating Correlation and plotting the Heat map of the correlated data  
important_correlation = important_df.corr()  
plt.figure(figsize=(25,25))  
sns.heatmap(important_correlation)  
plt.title("Correlation of the important features extracted")  
plt.show()
```



We see high correlation between few sensors with a value almost above 0.99 such as

- Sensor 1 and Sensor 45
 - Sensor 8 and Sensor 32
 - Sensor 27 and Sensor 46
 - Sensor 46 and Sensor 47
 - Sensor 47 and Sensor 67

We can drop those features that are highly correlated and keep only one among them

```
In [52]: for i, column_1 in enumerate(important_correlation.columns):
    for j in range(i, len(important_correlation.columns)):
        column_2 = important_correlation.columns[j]
        if column_1 != column_2:
            if important_correlation[column_1][column_2] > 0.99:
                print(column_1, "--", column_2, "Correlation Value: ", important_correlation[column_1][column_2])
```

sensor1_measure -- sensor45_measure Correlation Value: 0.9987741087714981
sensor8_measure -- sensor32_measure Correlation Value: 0.9977448921173756
sensor27_measure -- sensor46_measure Correlation Value: 0.9962581925040437
sensor27_measure -- sensor47_measure Correlation Value: 0.9962581925084124
sensor27_measure -- sensor67_measure Correlation Value: 0.996258192484215
sensor46_measure -- sensor47_measure Correlation Value: 0.9999999999896727
sensor46_measure -- sensor67_measure Correlation Value: 0.9999999999999973

Molecular mechanisms of the circadian clock in plants

```
In [53]: histogram_columns = dict()

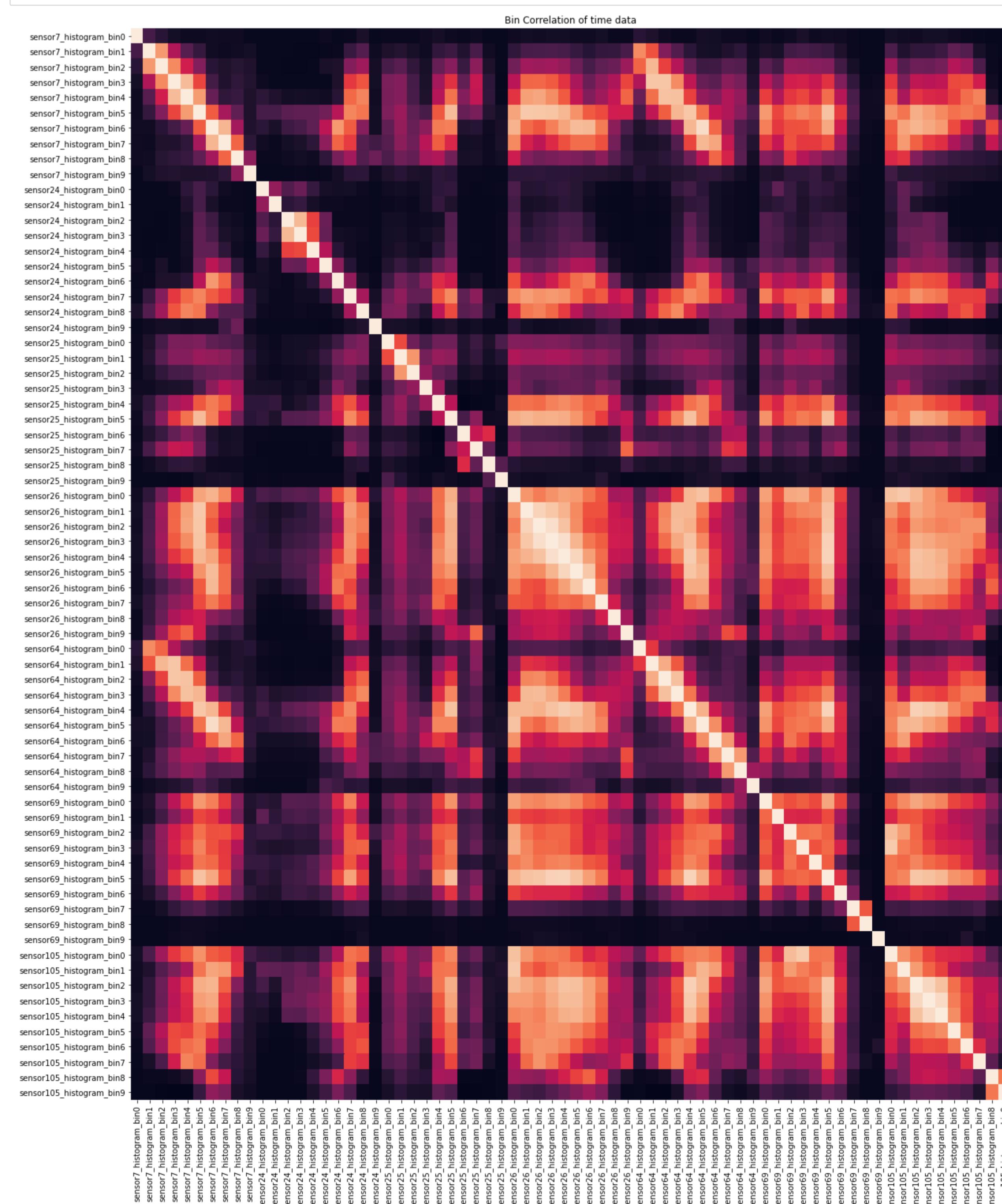
for column in data_df_copy.columns:
    value = column.split("_")
    if value[1] == 'histogram':
        if value[0] not in histogram_columns:
            histogram_columns[value[0]] = []
        histogram_columns[value[0]].append(column)
print(histogram_columns)

{'sensor7': ['sensor7_histogram_bin0', 'sensor7_histogram_bin1', 'sensor7_histogram_bin2', 'sensor7_histogram_bin3', 'sensor7_histogram_bin4', 'sensor7_histogram_bin5', 'sensor7_histogram_bin6', 'sensor7_histogram_bin7', 'sensor7_histogram_bin8', 'sensor7_histogram_bin9'], 'sensor24': ['sensor24_histogram_bin0', 'sensor24_histogram_bin1', 'sensor24_histogram_bin2', 'sensor24_histogram_bin3', 'sensor24_histogram_bin4', 'sensor24_histogram_bin5', 'sensor24_histogram_bin6', 'sensor24_histogram_bin7', 'sensor24_histogram_bin8', 'sensor24_histogram_bin9'], 'sensor25': ['sensor25_histogram_bin0', 'sensor25_histogram_bin1', 'sensor25_histogram_bin2', 'sensor25_histogram_bin3', 'sensor25_histogram_bin4', 'sensor25_histogram_bin5', 'sensor25_histogram_bin6', 'sensor25_histogram_bin7', 'sensor25_histogram_bin8', 'sensor25_histogram_bin9'], 'sensor26': ['sensor26_histogram_bin0', 'sensor26_histogram_bin1', 'sensor26_histogram_bin2', 'sensor26_histogram_bin3', 'sensor26_histogram_bin4', 'sensor26_histogram_bin5', 'sensor26_histogram_bin6', 'sensor26_histogram_bin7', 'sensor26_histogram_bin8', 'sensor26_histogram_bin9'], 'sensor64': ['sensor64_histogram_bin0', 'sensor64_histogram_bin1', 'sensor64_histogram_bin2', 'sensor64_histogram_bin3', 'sensor64_histogram_bin4', 'sensor64_histogram_bin5', 'sensor64_histogram_bin6', 'sensor64_histogram_bin7', 'sensor64_histogram_bin8', 'sensor64_histogram_bin9'], 'sensor69': ['sensor69_histogram_bin0', 'sensor69_histogram_bin1', 'sensor69_histogram_bin2', 'sensor69_histogram_bin3', 'sensor69_histogram_bin4', 'sensor69_histogram_bin5', 'sensor69_histogram_bin6', 'sensor69_histogram_bin7', 'sensor69_histogram_bin8', 'sensor69_histogram_bin9'], 'sensor105': ['sensor105_histogram_bin0', 'sensor105_histogram_bin1', 'sensor105_histogram_bin2', 'sensor105_histogram_bin3', 'sensor105_histogram_bin4', 'sensor105_histogram_bin5', 'sensor105_histogram_bin6', 'sensor105_histogram_bin7', 'sensor105_histogram_bin8', 'sensor105_histogram_bin9']}
```

```
In [54]: bins = []
for column, histograms in histogram_columns.items():
    bins.extend(histograms)
    i += 1
```

```
print(len(bins))  
70  
  
In [55]: histogram_data = data_df_copy[bins]  
print(histogram_data.shape)  
(60000, 70)  
  
In [56]: # Calculating Correlation and plotting the Heat map of the correlated data  
bin_correlation = histogram_data.corr()  
plt.figure(figsize=(25,25))
```

```
sns.heatmap(bin_correlation)
plt.title("Bin Correlation of time data")
plt.show()
```



- We see a correlation between the same sensor bins as they are continuous in times
- We also see high correlation between bins of different sensors such as Sensor 7 bin 6 and Sensor 64 bin 5 with correlation value above 0.95

Columns with high correlation data with some other column are as follows

```
In [57]: for i, column_1 in enumerate(bin_correlation.columns):
    for j in range(i, len(bin_correlation.columns)):
        column_2 = bin_correlation.columns[j]
        if column_1 != column_2:
            if bin_correlation[column_1][column_2] > 0.95:
                print(column_1, "--", column_2, "Correlation Value: ", bin_correlation[column_1][column_2])

sensor7_histogram_bin6 -- sensor64_histogram_bin5 Correlation Value: 0.9567868937028872
sensor26_histogram_bin2 -- sensor26_histogram_bin4 Correlation Value: 0.9649715653249844
sensor26_histogram_bin3 -- sensor26_histogram_bin4 Correlation Value: 0.96088366464621
sensor105_histogram_Bin3 -- sensor105_histogram_Bin4 Correlation Value: 0.9558991915181737
```

From observing the values, most bins which are next in time are the ones with high correlation since the sensor values are continuation of the previous bin values like "sensor7_histogram_bin6 -- sensor64_histogram_bin5".

The other 3 are correlation between subsequent bins of the same sensor due their time based relations

```
In [58]: for i, column_1 in enumerate(bin_correlation.columns):
    for j in range(i, len(bin_correlation.columns)):
        column_2 = bin_correlation.columns[j]
        if column_1 == column_2:
            sensor1 = column_1.split("_")[0]
            sensor2 = column_2.split("_")[0]
            if sensor1 == sensor2:
                if bin_correlation[column_1][column_2] > 0.95:
                    print(column_1, "--", column_2, "Correlation Value: ", bin_correlation[column_1][column_2])

sensor7_histogram_bin6 -- sensor64_histogram_bin5 Correlation Value: 0.9567868937028872
```

Data Preprocessing

```
In [59]: y_true = data_df['target']
data_df = data_df.drop(['target'], axis=1)
data_df = data_df.drop(['id'], axis=1)

In [60]: data_df.shape
Out[60]: (60000, 149)

In [61]: X_train, X_test, y_train, y_test = train_test_split(data_df, y_true, stratify=y_true, test_size=0.2)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)

In [62]: print(X_train.shape, X_cv.shape, X_test.shape)
(38400, 149) (9600, 149) (12000, 149)

From the above correlation values we see many sensors are correlated. We will remove the ones which are above 0.99 which means it is almost same as the other sensor and keep only among them.
• sensor1_measure : sensor45_measure Correlation Value: 0.9987741087714981
• sensor2d_measure : sensor32_measure Correlation Value: 0.99744892117375
• sensor27_measure : sensor46_measure Correlation Value: 0.9962581925040437
• sensor27_measure : sensor47_measure Correlation Value: 0.9962581925040437
• sensor27_measure : sensor67_measure Correlation Value: 0.9998999999999997
• sensor6d_measure : sensor67_measure Correlation Value: 0.9998999999999997
• sensor67_measure : sensor67d_measure Correlation Value: 0.999999999999997332
• sensor47_measure : sensor67d_measure Correlation Value: 0.999999999999997332
```

```
In [63]: high_correlated_columns = ['sensor45_measure', 'sensor32_measure', 'sensor46_measure', 'sensor47_measure', 'sensor67_measure']
for col in high_correlated_columns:
    if col in X_train.columns:
        X_train = X_train.drop([col], axis=1)
        X_cv = X_cv.drop([col], axis=1)
        X_test = X_test.drop([col], axis=1)
```

As discussed above, from the histogram bins, if there is a correlation between two different sensor bins then we should delete it. If there is correlation between the two bins of same sensor then it is due to subsequent time data.

```
• sensor7_histogram_bin6 : sensor64_histogram_bin5 Correlation Value: 0.9567868937028872
```

```
In [64]: X_train = X_train.drop(['sensor64_histogram_bin5'], axis=1)
X_cv = X_cv.drop(['sensor64_histogram_bin5'], axis=1)
X_test = X_test.drop(['sensor64_histogram_bin5'], axis=1)
```

After removing all the unwanted columns, we are left with 143 columns.

```
In [65]: print(X_train.shape, X_cv.shape, X_test.shape)
(38400, 143) (9600, 143) (12000, 143)
```

We have already dropped columns with large amount of zeros and NaN. Now we can impute the remaining NaN values with the median imputation.

Median imputation is performed since median is less error prone in this scenario where there is large fluctuation of sensor readings. Mean would provide a volatile value since either the values are zero or very large.

```
In [66]: median_values = X_train.median()
X_train_median = X_train.fillna(median_values)
X_cv_median = X_cv.fillna(median_values)
X_test_median = X_test.fillna(median_values)

print(X_train_median.shape, X_cv_median.shape, X_test_median.shape)
(38400, 143) (9600, 143) (12000, 143)

In [67]: dump(median_values, 'median_values.joblib')

Out[67]: ['median_values.joblib']

In [68]: with open('median_values.pkl', 'wb') as file:
pickle.dump(median_values, file)

In [69]: column_names = X_train.columns

In [69]: X_train_median
Out[69]:
```

	sensor1_measure	sensor3d_measure	sensor4_measure	sensor5_measure	sensor6_measure	sensor7_histogram_bin2	sensor7_histogram_bin3	sensor7_histogram_bin4	sensor7_histogram_bin5	sensor7_histogram_bin6	sensor105_histogram_bin1	sensor105_histogram_bin2	sensor105_histogram_bin3	sensor105_histogram_bin4	sensor105_histogram_bin5	sensor105_histogram_bin6	sensor105_histogram_bin7	sensor105_hist
40583	2256	3.000000e+01	20.0	0.0	0.0	0.0	3746.0	17810.0	93112.0	... 348980.0	234915.0	112901.0	221924.0	190452.0	94335.0	41190.0		
28732	1384	3.400000e+01	32.0	0.0	0.0	0.0	196.0	32770.0	45914.0	... 9196.0	4110.0	1048.0	9880.0	42716.0	4594.0	26.0		
15593	40854	1.520000e+02	80.0	0.0	0.0	0.0	7710.0	56124.0	1888174.0	... 431328.0	373538.0	178972.0	347670.0	312232.0	451418.0	146280.0		
34808	14	1.400000e+01	6.0	0.0	0.0	0.0	70.0	3730.0	2494.0	... 1740.0	106.0	56.0	68.0	50.0	60.0	48.0		
40025	28800	2.520000e+02	244.0	0.0	0.0	0.0	1876.0	20971.0	1291356.0	... 343080.0	309968.0	195206.0	402762.0	222626.0	94058.0	53634.0		
...		
42629	2236	2.400000e+01	126.0	0.0	0.0	0.0	38136.0	120410.0	9660.0	... 1262.0	21288.0	4676.0	2836.0	16214.0	26248.0	21202.0	23812.0	
49812	229694	1.500000e+02	126.0	0.0	0.0	0.0	525914.0	9314812.0	23660526.0	21448308.0	10097812.0	6529190.0	1968150.0	2020942.0	726960.0	307820.0	194850.0	
45178	101630	4.280000e+02	412.0	0.0	0.0	0.0	0.0	28276.0	2725512.0	3747300.0	1305642.0	721862.0	342108.0	680648.0	494300.0	459820.0	371714.0	
16101	40226	2.130706e+09	600.0	0.0	0.0	0.0	508.0	191428.0	2493750.0	597714.0	540596.0	401694.0	766502.0	215678.0	20880.0	16402.0		
53966	1704	2.800000e+01	18.0	0.0	0.0	0.0	0.0	20398.0	75480.0	4236.0	8124.0	3242.0	2230.0	61642.0	4538.0	1790.0	7252.0	

38400 rows x 143 columns

KNN based nearest neighbour imputation so as to provide more accurate imputation over simple imputations of median, mean, most frequent.

```
In [237]: knn_imputer = KNNImputer(weights="distance")
X_train_knn_imputation = knn_imputer.fit_transform(X_train)
X_cv_knn_imputation = knn_imputer.transform(X_cv)
X_test_knn_imputation = knn_imputer.transform(X_test)

print(X_train_knn_imputation.shape, X_cv_knn_imputation.shape, X_test_knn_imputation.shape)
(38400, 143) (9600, 143) (12000, 143)

In [238]: print(X_train_knn_imputation.shape, X_cv_knn_imputation.shape, X_test_knn_imputation.shape)
(38400, 143) (9600, 143) (12000, 143)

In [239]: X_train_knn_imputation = pd.DataFrame(X_train_knn_imputation)
X_cv_knn_imputation = pd.DataFrame(X_cv_knn_imputation)
X_test_knn_imputation = pd.DataFrame(X_test_knn_imputation)

X_train_knn_imputation.columns = column_names
X_cv_knn_imputation.columns = column_names
X_test_knn_imputation.columns = column_names
```

Feature Engineering

Addition of 2 columns using Truncated SVD as part of feature engineering for median and knn imputed values.

```
In [71]: T_SVD = TruncatedSVD(n_components=4, n_iter=20)
X_train_median_tSVD = T_SVD.fit_transform(X_train_median)
X_cv_median_tSVD = T_SVD.transform(X_cv_median)
X_test_median_tSVD = T_SVD.fit_transform(X_test_median)

print(X_train_median_tSVD.shape, X_cv_median_tSVD.shape, X_test_median_tSVD.shape)
(38400, 4) (9600, 4) (12000, 4)

In [72]: dump(T_SVD, 'truncated_SVD.joblib')
Out[72]: ['truncated_SVD.joblib']

In [90]: with open('truncated_SVD.pkl', 'wb') as file:
pickle.dump(T_SVD, file)

In [73]: for i in range(4):
    X_train_median["SVD_"+str(i)] = X_train_median_tSVD[:,i]
    X_cv_median["SVD_"+str(i)] = X_cv_median_tSVD[:,i]
    X_test_median["SVD_"+str(i)] = X_test_median_tSVD[:,i]

print(X_train_median.shape, X_cv_median.shape, X_test_median.shape)
(38400, 147) (9600, 147) (12000, 147)

In [242]: T_SVD = TruncatedSVD(n_components=4, n_iter=20)
X_train_knn_imputation_tSVD = T_SVD.fit_transform(X_train_knn_imputation)
X_cv_knn_imputation_tSVD = T_SVD.transform(X_cv_knn_imputation)
X_test_knn_imputation_tSVD = T_SVD.fit_transform(X_test_knn_imputation)

print(X_train_knn_imputation_tSVD.shape, X_cv_knn_imputation_tSVD.shape, X_test_knn_imputation_tSVD.shape)
(38400, 4) (9600, 4) (12000, 4)

In [243]: for i in range(4):
    X_train_knn_imputation["SVD_"+str(i)] = X_train_knn_imputation_tSVD[:,i]
    X_cv_knn_imputation["SVD_"+str(i)] = X_cv_knn_imputation_tSVD[:,i]
    X_test_knn_imputation["SVD_"+str(i)] = X_test_knn_imputation_tSVD[:,i]

print(X_train_knn_imputation.shape, X_cv_knn_imputation.shape, X_test_knn_imputation.shape)
(38400, 147) (9600, 147) (12000, 147)
```

Add bin average as a separate column

Non ML based Modelling

Currently there are

Random Model

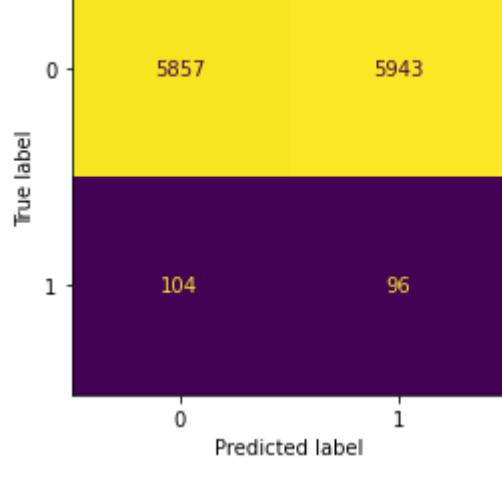
```
In [263]: y_train_pred, y_cv_pred, y_test_pred = random_model(X_train_knn_os.shape[0], X_cv_knn_imputation.shape[0], X_test_knn_imputation.shape[0])

print("Random Model: Train F1 score:",f1_score(y_train_knn_os, y_train_pred, average='macro'))
print("Random Model: CV F1 score:",f1_score(y_cv, y_cv_pred, average='macro'))
print("Random Model: Test F1 score:",f1_score(y_test, y_test_pred, average='macro'))

cm = confusion_matrix(y_test, y_test_pred, labels=[0,1])
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0,1])
disp.plot()

Random Model: Train F1 score: 0.3994899335055624
Random Model: CV F1 score: 0.3488472242359573
Random Model: Test F1 score: 0.34515454931099315
```

Out[263]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x222b1137a48>



Predicting All 1s

```
In [264]: # Predict all the values as 1s
def model_ones(train_shape, cv_shape, test_shape):
    train = np.ones(train_shape)
    cv = np.ones(cv_shape)
    test = np.ones(test_shape)
    return train, cv, test

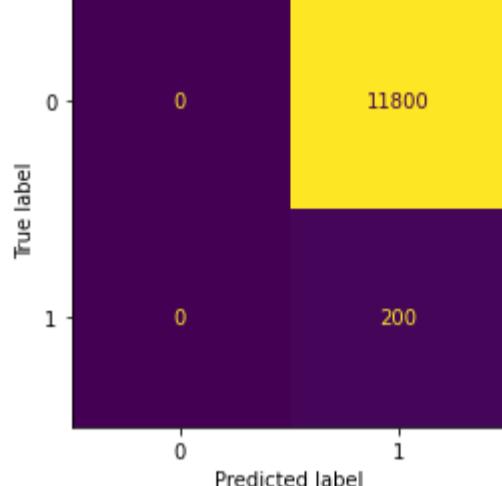
y_train_pred, y_cv_pred, y_test_pred = model_ones(X_train_knn_os.shape[0], X_cv_knn_imputation.shape[0], X_test_knn_imputation.shape[0])

print("All 1s Model: Train F1 score:",f1_score(y_train_knn_os, y_train_pred, average='macro'))
print("All 1s Model: CV F1 score:",f1_score(y_cv, y_cv_pred, average='macro'))
print("All 1s Model: Test F1 score:",f1_score(y_test, y_test_pred, average='macro'))

cm = confusion_matrix(y_test, y_test_pred, labels=[0,1])
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0,1])
disp.plot()

All 1s Model: Train F1 score: 0.9833333333333334
All 1s Model: CV F1 score: 0.9139344262295082
All 1s Model: Test F1 score: 0.91639344262295082
```

Out[264]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x222b1bfdf88>



Model predicting all 1s will result in a bad f1 score because the class with 1s is significantly low and hence the precision will be highly low and thus resulting in a low f1 score.

Predicting all 0s

```
In [265]: # Predict all the values as Zeros
def model_zeros(train_shape, cv_shape, test_shape):
    train = np.zeros(train_shape)
    cv = np.zeros(cv_shape)
    test = np.zeros(test_shape)
    return train, cv, test

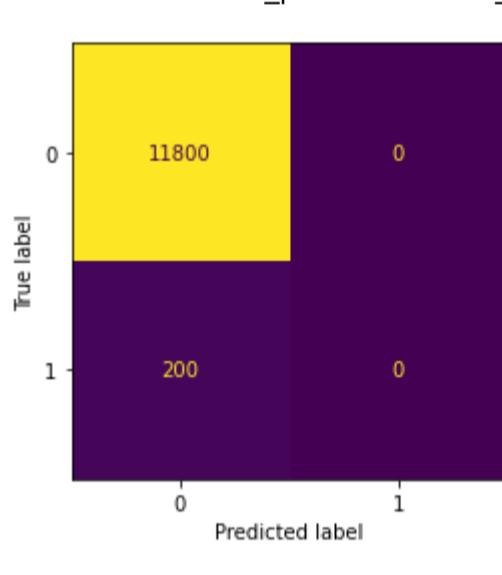
y_train_pred, y_cv_pred, y_test_pred = model_zeros(X_train_knn_os.shape[0], X_cv_knn_imputation.shape[0], X_test_knn_imputation.shape[0])

print("All 0s Model: Train F1 score:",f1_score(y_train_knn_os, y_train_pred, average='macro'))
print("All 0s Model: CV F1 score:",f1_score(y_cv, y_cv_pred, average='macro'))
print("All 0s Model: Test F1 score:",f1_score(y_test, y_test_pred, average='macro'))

cm = confusion_matrix(y_test, y_test_pred, labels=[0,1])
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0,1])
disp.plot()

All 0s Model: Train F1 score: 0.47619047619047616
All 0s Model: CV F1 score: 0.49579831932773105
All 0s Model: Test F1 score: 0.49579831932773105
```

Out[265]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x222b1fa00c8>



When predicting all zeros we tend to get a f1 score of ~0.5 this is to the fact that the entire minority class is getting misclassified.

```
In [266]: X_train_knn_os.head()
Out[266]: sensor1_measure sensor3_measure sensor4_measure sensor5_measure sensor6_measure sensor7_histogram_bin2 sensor7_histogram_bin3 sensor7_histogram_bin4 sensor7_histogram_bin5 sensor7_histogram_bin6 ... SVD_1 SVD_2 SVD_3 sensor7_bin_average sensor24_bin_average sensor25_bin_average sensor64_bin_average sensor69_bin_average sensor105_bin_average
0 0.01432 0.268456e-08 0.127467e-08 0.0 0.0 0.00000 0.00000 0.000036 0.007690 0.019205 ... 0.004513 0.000661 0.282686 0.006111 0.006113 0.006111 0.005361 0.006111 0.006111
1 0.00031 0.00000e+00 0.00000e+00 0.0 0.0 0.00000 0.00000 0.00000 0.000078 0.000045 ... 0.004272 0.00061 0.280331 0.000035 0.000035 0.000035 0.000041 0.000035 0.000035
2 0.00056 1.595715e-08 3.379160e-08 0.745547e-09 0.0 0.0 0.00000 0.00000 0.000021 0.000275 0.000024 0.000024 0.000015 0.000015 ... 0.004271 0.00061 0.28027 0.000048 0.000048 0.000048 0.000052 0.000048 0.000048
3 0.01490 3.379160e-08 0.745547e-09 0.0 0.0 0.00000 0.00000 0.00000 0.000767 0.000676 0.010983 ... 0.004455 0.00061 0.282401 0.004108 0.004108 0.004108 0.003873 0.004108 0.004108
4 0.01453 9.999986e-01 2.819101e-08 0.0 0.0 0.00000 0.00000 0.00000 0.000006 0.003011 0.012316 ... 0.000189 0.00061 0.274487 0.004128 0.004128 0.004128 0.003274 0.004128 0.004128
```

5 rows x 154 columns

Model based on percentiles of certain sensors where classes can be separated based on percentile values

```
In [267]: sensor_names = ['sensor1_measure', 'sensor8_measure', 'sensor14_measure', 'sensor15_measure', 'sensor16_measure', 'sensor17_measure', 'sensor27_measure', 'sensor33_measure', 'sensor34_measure', 'sensor35_measure', 'sensor59_measure', 'sensor61_measure', 'sensor89_measure']

def sensor_based_prediction(X_train, X_cv, X_test, y_true, y_cv, y_test):

    train_f1_scores = []
    cv_f1_scores = []
    test_f1_scores = []
    final_test_scor = dict()

    X_train['label'] = y_true

    df_0 = X_train[X_train['label']==0]
    df_1 = X_train[X_train['label']==1]

    df_0 = df_0[sensor_names]
    df_1 = df_1[sensor_names]

    quantiles_0 = df_0.quantile([0.1, 0.9], axis=0)
    quantiles_1 = df_1.quantile([0.1, 0.9], axis=0)

    # For each sensor column
    for sensor_name in sensor_names:
        y_train_pred = []
        y_cv_pred = []
        y_test_pred = []

        # If the value is less than the (90th quantile + 10th 1 quantile)/2 then set the value as class 0 else it is class 1
        for value in X_train[sensor_name]:
            if value < ((quantiles_0[sensor_name][0.9] + quantiles_1[sensor_name][0.1])/2):
                y_train_pred.append(0)
            else:
                y_train_pred.append(1)

        # Loop same for CV dataset
        for value in X_cv[sensor_name]:
            if value < ((quantiles_0[sensor_name][0.9] + quantiles_1[sensor_name][0.1])/2):
                y_cv_pred.append(0)
            else:
                y_cv_pred.append(1)

        # Loop same for test dataset
        for value in X_test[sensor_name]:
            if value < ((quantiles_0[sensor_name][0.9] + quantiles_1[sensor_name][0.1])/2):
                y_test_pred.append(0)
            else:
                y_test_pred.append(1)

        print(sensor_name, ": f1 score : ", f1_score(y_true, np.array(y_train_pred), average='macro'), \
              f1_score(y_cv, np.array(y_cv_pred), average='macro'), \
              f1_score(y_test, np.array(y_test_pred), average='macro'))

        final_test_scor[sensor_name] = f1_score(y_test, np.array(y_test_pred), average='macro')

    return final_test_scor
```

For KNN based imputation

```
In [268]: f1_scores_knn_os = sensor_based_prediction(X_train_knn_os, X_cv_knn_imputation, X_test_knn_imputation, y_train_knn_os, y_cv, y_test)
```

```
sensor1_measure : f1 score : 0.8392836994940816 0.6481773238651789 0.63762345399287
sensor2_measure : f1 score : 0.839346409920944 0.633895729132079 0.64034468716934
sensor3_measure : f1 score : 0.83364935382083 0.6252553729138072 0.6349116763598877
sensor4_measure : f1 score : 0.8334671995170909 0.592732628230939 0.6834837165391876
sensor5_measure : f1 score : 0.845417995588668 0.6491017896244606 0.64827159939354153
sensor6_measure : f1 score : 0.8436413937536121 0.6307312162536159 0.64325438381616
sensor7_measure : f1 score : 0.8436413937536121 0.6307312162536159 0.64325438381616
sensor8_measure : f1 score : 0.8448999025593121 0.6453090161367632 0.6428390242979268
sensor9_measure : f1 score : 0.8125875295245186 0.6077035236159822 0.6123524514860155
sensor10_measure : f1 score : 0.856499535923799 0.659524926293219 0.65626678318809
sensor11_measure : f1 score : 0.844891851879279 0.6433407683167239 0.643374768239343
sensor12_measure : f1 score : 0.847115882468599 0.6510793148728815 0.645792555265859
sensor13_measure : f1 score : 0.8576439439210823 0.655305919844681 0.649388395696175
```

We see that Sensor 35, Sensor 61 and Sensor 89 are the top 3 f1 scores. Thus these may seem to be more important in deciding the class when a ML based Model is used

For Median based imputation

```
In [269]: f1_scores_median_os = sensor_based_prediction(X_train_median_os, X_cv_median, X_test_median, y_train_median_os, y_cv, y_test)
```

```
sensor1_measure : f1 score : 0.8382836994940816 0.6481773238651789 0.63762345399287
sensor2_measure : f1 score : 0.839346409920944 0.633895729132079 0.64034468716934
sensor3_measure : f1 score : 0.79339670545071994 0.592732628230939 0.6834837165391876
sensor4_measure : f1 score : 0.83364935382083 0.6252553729138072 0.6349116763598877
sensor5_measure : f1 score : 0.8152978741177985 0.617094358392558 0.623544698157617
sensor6_measure : f1 score : 0.8081199605122134 0.607787595193219 0.611877129965237
sensor7_measure : f1 score : 0.8048228617377984 0.604359087694945 0.613643767032781
sensor8_measure : f1 score : 0.811718316953812 0.6096877038838664 0.61728844998019587
sensor9_measure : f1 score : 0.8383369936918634 0.6312435923231843 0.635887719292456
sensor10_measure : f1 score : 0.843226448822774 0.6432880289776763 0.6437002123499
sensor11_measure : f1 score : 0.848914731950244 0.6479770364297448 0.6458119541657708
sensor12_measure : f1 score : 0.8270920678874041 0.626207366996489 0.6282368345760841
```

For median based imputation, we see that Sensor 61, Sensor 59, Sensor 1 and Sensor 35 are the important sensor based on the F1 scores alone.

We can also conclude that the KNN based imputation provided a slightly better imputation, the same result can also be seen in their respective F1 score of each sensors.

```
In [270]: X_train_median_os = X_train_median_os.drop(['label'], axis=1)
X_train_knn_os = X_train_knn_os.drop(['label'], axis=1)
```

ML Modelling

KNN

Median imputed data

```
In [93]: # KNN with Random Search CV for hyper parameter tuning
x_cfl=KNeighborsClassifier()

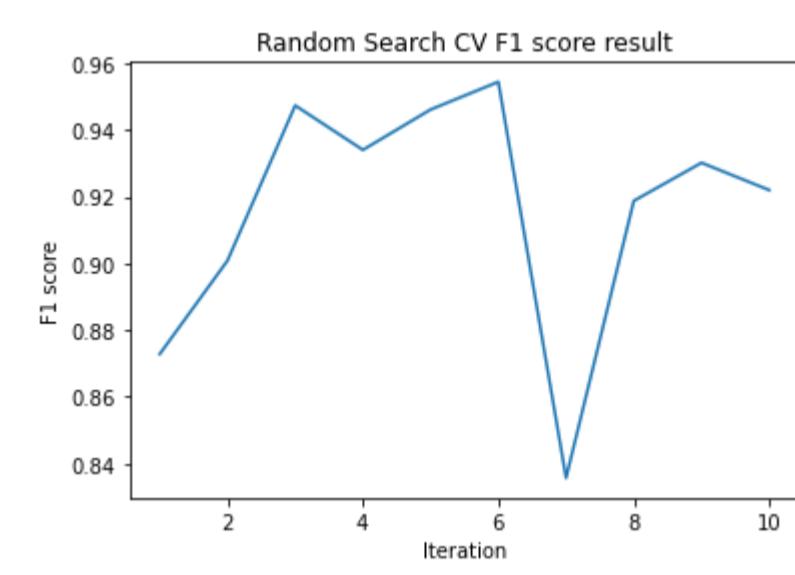
params={
    'n_neighbors':[3, 4, 10, 20, 50],
    'weights':['uniform', 'distance'],
    'metric':['euclidean', 'manhattan']
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=params,n_jobs=-1, cv=3, scoring = "f1")
random_cfl.fit(X_train_median_os,y_train_median_os.values.ravel())

Out[93]: RandomizedSearchCV(cv=3, estimator=KNeighborsClassifier(n_neighbors=3,
                                                               metric='euclidean', weights='uniform'),
                           param_distributions={'n_neighbors': [3, 4, 10, 20, 50],
                                                 'weights': ['uniform', 'distance']},
                           scoring='f1')
```

```
In [94]: # Displaying mean test score which is the mean F1 score for each hyper parameter combination
labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)

Out[94]: <matplotlib.lines.Line2D at 0x248b3be448>
```



```
In [95]: random_cfl.best_estimator_
Out[95]: KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')

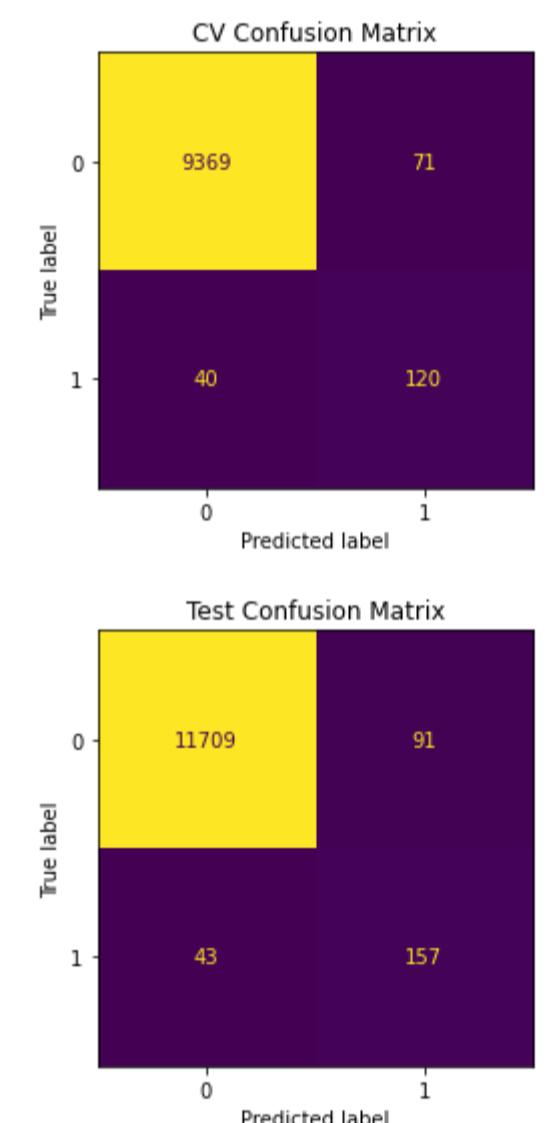
In [274]: x_cfl=KNeighborsClassifier(metric='manhattan', n_neighbors=3, weights='distance')
x_cfl.fit(X_train_median_os,y_train_median_os.values.ravel())

Out[274]: KNeighborsClassifier(metric='manhattan', n_neighbors=3, weights='distance')
```

```
In [275]: predict_y = x_cfl.predict(X_train_median_os)
print ("Train F1 score:", f1_score(y_train_median_os, predict_y, average='macro'))
predict_y = x_cfl.predict(X_cv_median, y_cv)
print ("CV F1 score:", f1_score(y_cv, predict_y, average='macro'))
predict_y = x_cfl.predict(X_test_median)
print("Test F1 score:", f1_score(y_test, predict_y, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.83935885936954
Test F1 score: 0.8476015595496895
```

```
In [276]: plot1 = plot_confusion_matrix(x_cfl, X_cv_median, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')
plot2 = plot_confusion_matrix(x_cfl, X_test_median, y_test)
plot2.ax_.set_title('Test Confusion Matrix')
plot1.show()
```



KNN imputed data

```
In [277]: # KNN with Random Search CV for hyper parameter tuning
x_cfl=KNeighborsClassifier()

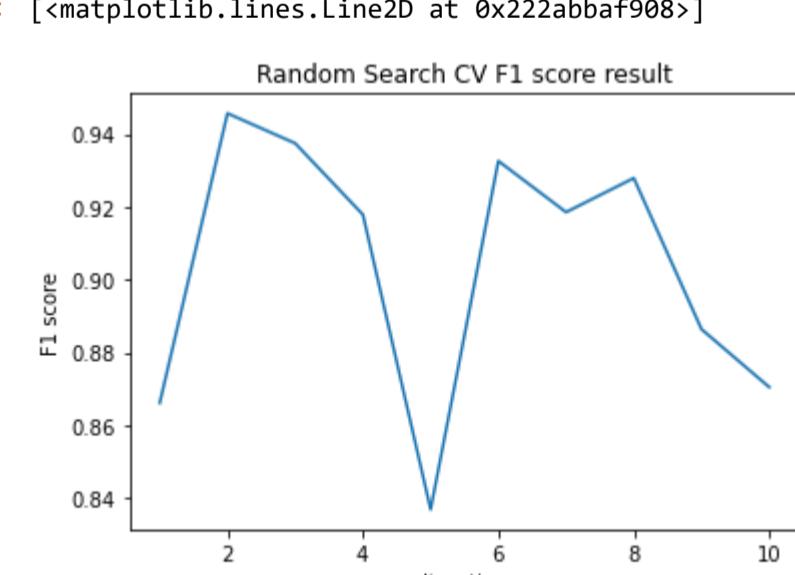
params={
    'n_neighbors':[3, 4, 10, 20, 50],
    'weights':['uniform', 'distance'],
    'metric':['euclidean', 'manhattan']
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=params,n_jobs=-1, cv=3, scoring = "f1")
random_cfl.fit(X_train_knn_os,y_train_knn_os.values.ravel())

Out[277]: RandomizedSearchCV(cv=3, estimator=KNeighborsClassifier(n_neighbors=4,
                                                               metric='euclidean', weights='uniform'),
                           param_distributions={'n_neighbors': [3, 4, 10, 20, 50],
                                                 'weights': ['uniform', 'distance']},
                           scoring='f1')
```

```
In [278]: # Displaying mean test score which is the mean F1 score for each hyper parameter combination
labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)

Out[278]: <matplotlib.lines.Line2D at 0x222abbaf908>
```



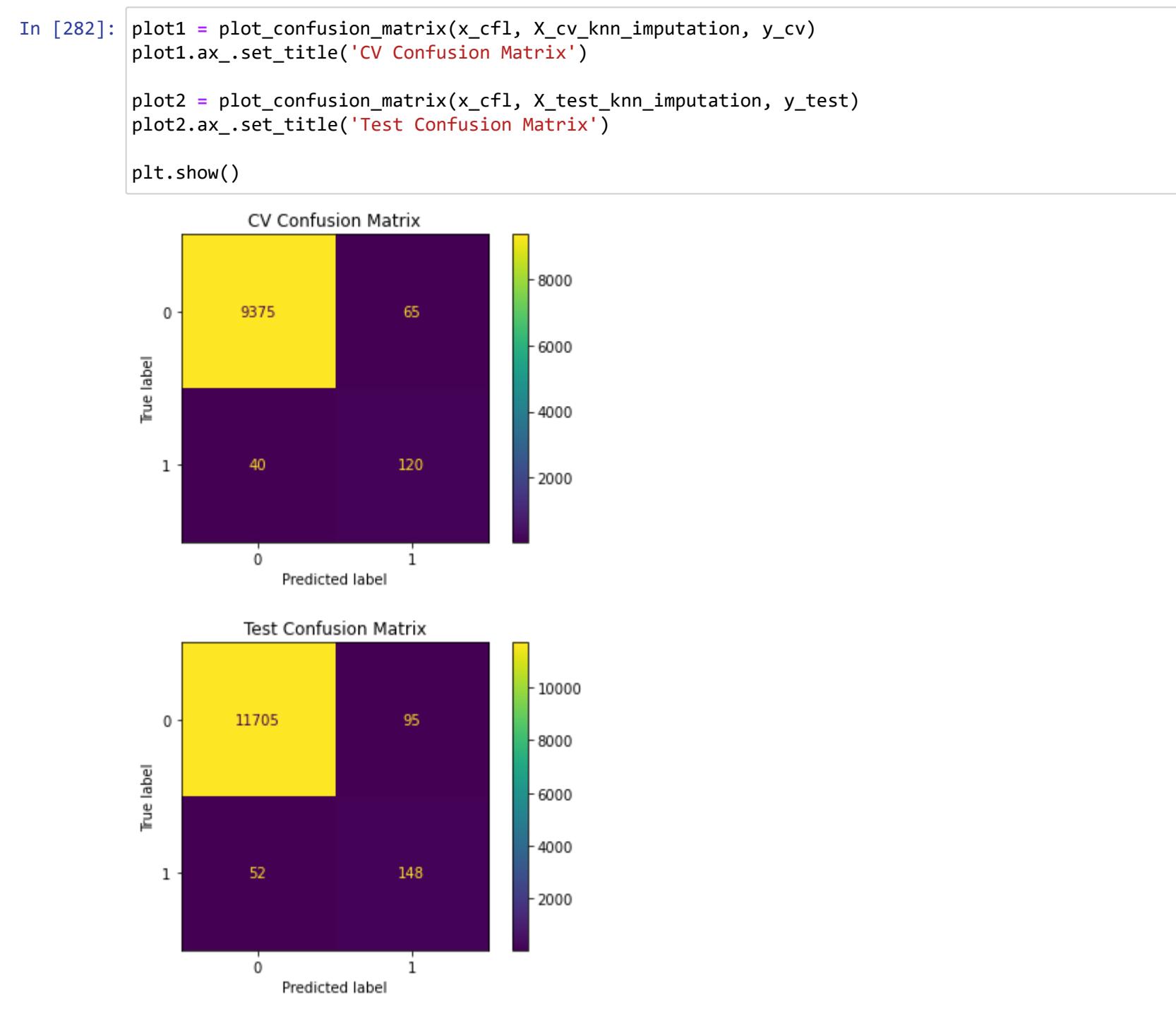
```
In [279]: random_cfl.best_estimator_
Out[279]: KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')

In [280]: x_cfl=KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')
x_cfl.fit(X_train_knn_os,y_train_knn_os.values.ravel())

Out[280]: KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')
```

```
In [281]: predict_y = x_cfl.predict(X_train_knn_os)
print ("Train F1 score:", f1_score(y_train_knn_imputation, predict_y, average='macro'))
predict_y = x_cfl.predict(X_cv_knn_imputation, y_cv)
print ("CV F1 score:", f1_score(y_cv, predict_y, average='macro'))
predict_y = x_cfl.predict(X_test_knn_imputation)
print("Test F1 score:", f1_score(y_test, predict_y, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.8450416796375982
Test F1 score: 0.8309656870885478
```



For KNN, We see that Median based imputations provided slightly better results compared to the KNN based imputation dataset.

However for the minority class, both the median based imputation dataset provided better classification.

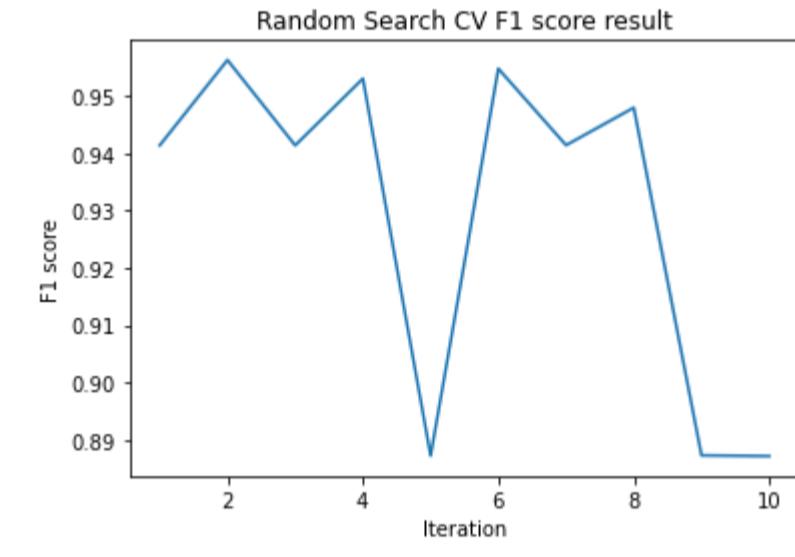
Random Forest Classifier

Median Imputed Data

```
In [96]: # Random Forest Classifier with Random Search CV for hyper parameter tuning
x_cfl=RandomForestClassifier()
params={
    "n_estimators": [100, 200, 400, 800, 1000, 1500, 2000],
    "max_depth": [5, 10, 20, 40, 100, 200],
    "min_samples_split": [2, 5, 10],
}
random_cfl = RandomizedSearchCV(x_cfl,param_distributions=params,n_jobs=-1, cv=3, scoring = "f1")
random_cfl.fit(X_train_median_os,y_train_median_os.values.ravel())
Out[96]: RandomizedSearchCV(cv=3, estimator=RandomForestClassifier(), n_jobs=-1,
                           param_distributions={"max_depth": [5, 10, 20, 40, 100, 200],
                                                 "min_samples_split": [2, 5, 10],
                                                 "n_estimators": [100, 200, 400, 800,
                                                 1000, 1500, 2000]},
                           scoring='f1')

In [97]: # Displaying mean test score which is the mean F1 score for each hyper parameter combination
labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])
    plt.xlabel('Iteration')
    plt.ylabel('F1 score')
    plt.title('Random Search CV F1 score result')
    plt.plot(x_axis, y_axis)

{'n_estimators': 1500, 'min_samples_split': 2, 'max_depth': 10} : 0.9413489939937195
{'n_estimators': 2000, 'min_samples_split': 5, 'max_depth': 10} : 0.93621011797172
{'n_estimators': 400, 'min_samples_split': 5, 'max_depth': 10} : 0.9413424463493959
{'n_estimators': 1500, 'min_samples_split': 5, 'max_depth': 40} : 0.9529810175765885
{'n_estimators': 800, 'min_samples_split': 5, 'max_depth': 5} : 0.8872108969284476
{'n_estimators': 200, 'min_samples_split': 2, 'max_depth': 200} : 0.954722729825244
{'n_estimators': 1000, 'min_samples_split': 5, 'max_depth': 10} : 0.9479218234922366
{'n_estimators': 200, 'min_samples_split': 10, 'max_depth': 30} : 0.9479218234922366
{'n_estimators': 2000, 'min_samples_split': 2, 'max_depth': 5} : 0.887278347362721
{'n_estimators': 1000, 'min_samples_split': 10, 'max_depth': 5} : 0.8871448967638166
Out[97]: <matplotlib.lines.Line2D at 0x248c98ec908>
```



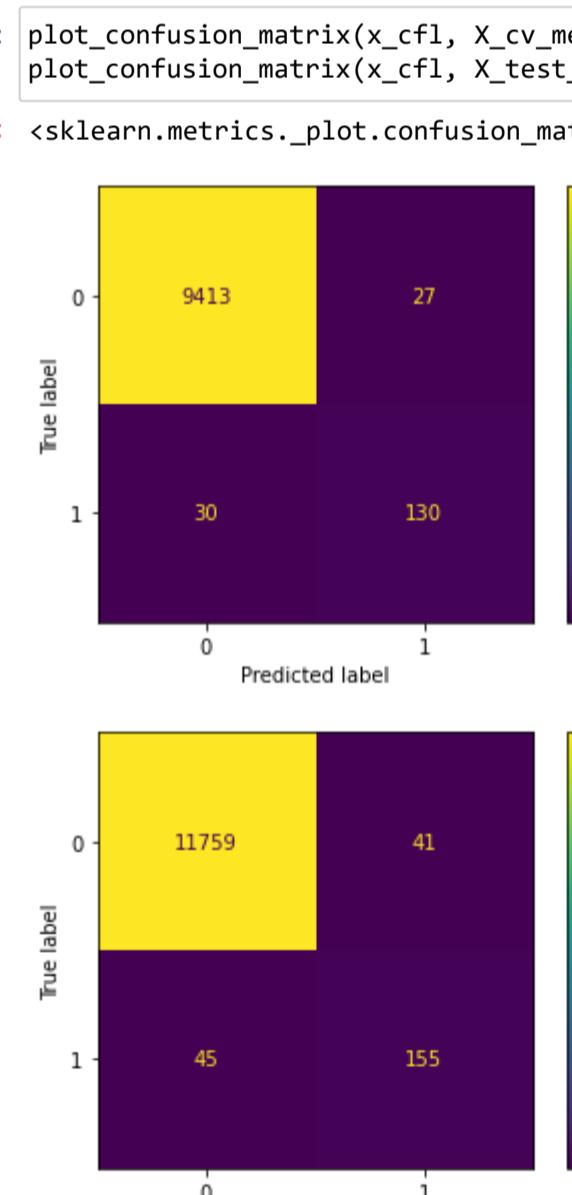
```
In [98]: random_cfl.best_estimator_
Out[98]: RandomForestClassifier(max_depth=100, n_estimators=2000)

In [206]: X_cfl = RandomForestClassifier(max_depth=100, n_estimators=1500)
X_cfl.fit(X_train_median_os,y_train_median_os.values.ravel())
Out[206]: RandomForestClassifier(max_depth=100, n_estimators=1500)
```

```
In [287]: predict_y = X_cfl.predict(X_train_median_os)
print("Train F1 score: ", f1_score(y_train,y_train_median))
predict_y = X_cfl.predict(X_cv.median)
print("CV F1 score: ", f1_score(y_cv, predict_y, average="macro"))
predict_y = X_cfl.predict(X_test_median)
print("Test F1 score: ", f1_score(y_test, predict_y, average="macro"))

Train F1 score: 1.0
CV F1 score: 0.908583431499399
Test F1 score: 0.8895924162828078
```

```
In [288]: plot_confusion_matrix(x_cfl, X_cv_median, y_cv)
plot_confusion_matrix(x_cfl, X_test_median, y_test)
Out[288]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x22a6c09e88>
```

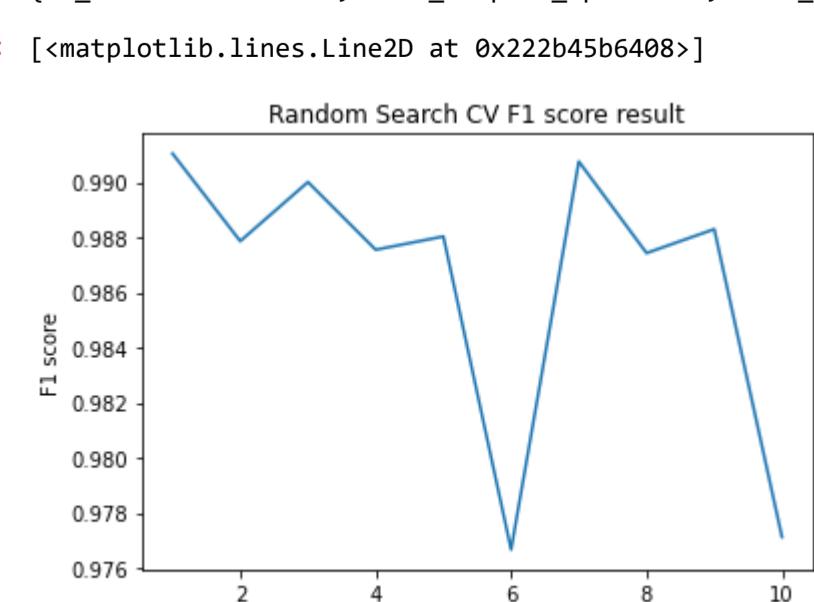


KNN Imputed data

```
In [289]: # Random Forest Classifier with Random Search CV for hyper parameter tuning
x_cfl = RandomForestClassifier()
params={
    "n_estimators": [100, 200, 400, 800, 1000, 1500, 2000],
    "max_depth": [5, 10, 20, 40, 100, 200],
    "min_samples_split": [2, 5, 10],
}
random_cfl = RandomizedSearchCV(x_cfl,param_distributions=params,verbose=10,n_jobs=-1, cv=3)
random_cfl.fit(X_train_knn_os,y_train_knn_os.values.ravel())
Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=1): Using backend LokyBackend with 12 concurrent workers.
[Parallel(n_jobs=1): Done 1 tasks | elapsed: 48.2s
[Parallel(n_jobs=1): Done 11 out of 30 | elapsed: 5.9min remaining: 10.3min
[Parallel(n_jobs=1): Done 15 out of 30 | elapsed: 6.3min remaining: 6.3min
[Parallel(n_jobs=1): Done 19 out of 30 | elapsed: 6.8min remaining: 3.9min
[Parallel(n_jobs=1): Done 23 out of 30 | elapsed: 7.3min remaining: 2.4min
[Parallel(n_jobs=1): Done 27 out of 30 | elapsed: 8.3min remaining: 54.9s
[Parallel(n_jobs=1): Done 30 out of 30 | elapsed: 10.2min finished
Out[289]: RandomizedSearchCV(cv=3, estimator=RandomForestClassifier(), n_jobs=-1,
                           param_distributions={"max_depth": [5, 10, 20, 40, 100, 200],
                                                 "min_samples_split": [2, 5, 10],
                                                 "n_estimators": [100, 200, 400, 800,
                                                 1000, 1500, 2000]},
                           verbose=10)
```

```
In [290]: labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])
    plt.xlabel('Iteration')
    plt.ylabel('F1 score')
    plt.title('Random Search CV F1 score result')
    plt.plot(x_axis, y_axis)

{'n_estimators': 2000, 'min_samples_split': 2, 'max_depth': 100} : 0.99143953696512
{'n_estimators': 200, 'min_samples_split': 10, 'max_depth': 10} : 0.9878659875911332
{'n_estimators': 1500, 'min_samples_split': 10, 'max_depth': 100} : 0.99000870484372
{'n_estimators': 1500, 'min_samples_split': 10, 'max_depth': 10} : 0.98753638484397
{'n_estimators': 200, 'min_samples_split': 10, 'max_depth': 100} : 0.99000870484372
{'n_estimators': 1500, 'min_samples_split': 10, 'max_depth': 5} : 0.9766788762422173
{'n_estimators': 200, 'min_samples_split': 5, 'max_depth': 40} : 0.990755092394742
{'n_estimators': 100, 'min_samples_split': 10, 'max_depth': 10} : 0.98743263909098475
{'n_estimators': 1500, 'min_samples_split': 2, 'max_depth': 10} : 0.988293379112662
{'n_estimators': 100, 'min_samples_split': 5, 'max_depth': 5} : 0.9771283095981539
Out[290]: <matplotlib.lines.Line2D at 0x2224545b64088>
```



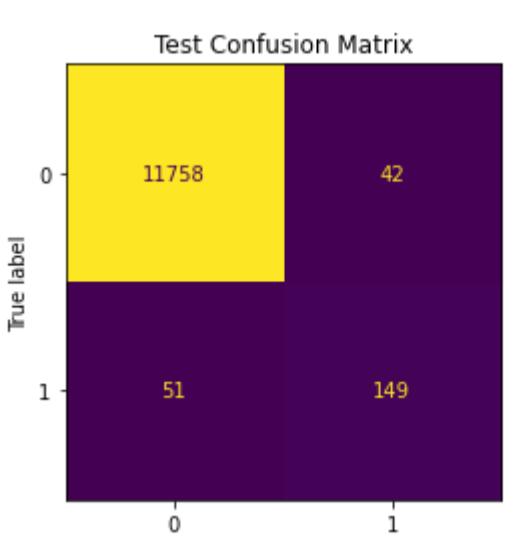
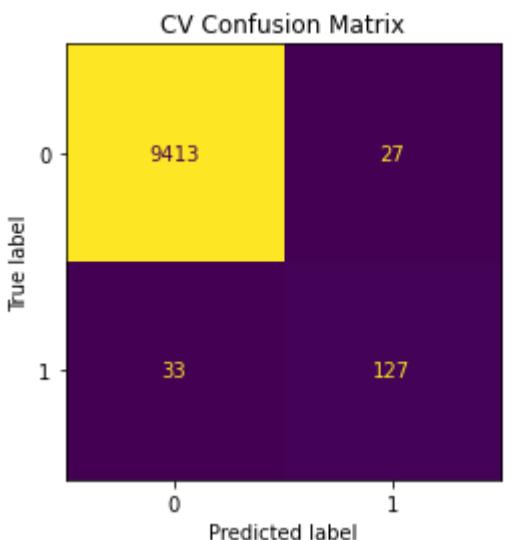
```
In [291]: random_cfl.best_estimator_
Out[291]: RandomForestClassifier(max_depth=100, n_estimators=1000)
```

```
In [292]: x_cfl = RandomForestClassifier(max_depth=100, n_estimators=1000)
x_cfl.fit(X_train_knn_os)
Out[292]: RandomForestClassifier(max_depth=100, n_estimators=1000)

In [293]: predict_y = x_cfl.predict(X_train_knn_os)
print("Train F1 score:", f1_score(y_train_knn_os, predict_y, average='macro'))
predict_y = x_cfl.predict(X_cv_knn_imputation)
print("CV F1 score:", f1_score(y_cv, predict_y, average='macro'))
predict_y = x_cfl.predict(X_test_knn_imputation)
print("Test F1 score:", f1_score(y_test, predict_y, average='macro'))

Train F1 score: 1.0
CV F1 score: 0.982870120482665
Test F1 score: 0.879104580928921

In [294]: plot1 = plot_confusion_matrix(x_cfl, X_cv_knn_imputation, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')
plot2 = plot_confusion_matrix(x_cfl, X_test_knn_imputation, y_test)
plot2.ax_.set_title('Test Confusion Matrix')
plot1.show()
```



The RF classifier with Median based imputation provided slightly better classification over the KNN based imputation dataset.

From KNN Models, RF classifier overall performed better overall but for the minority class, KNN Model had better performance and lesser misclassification.

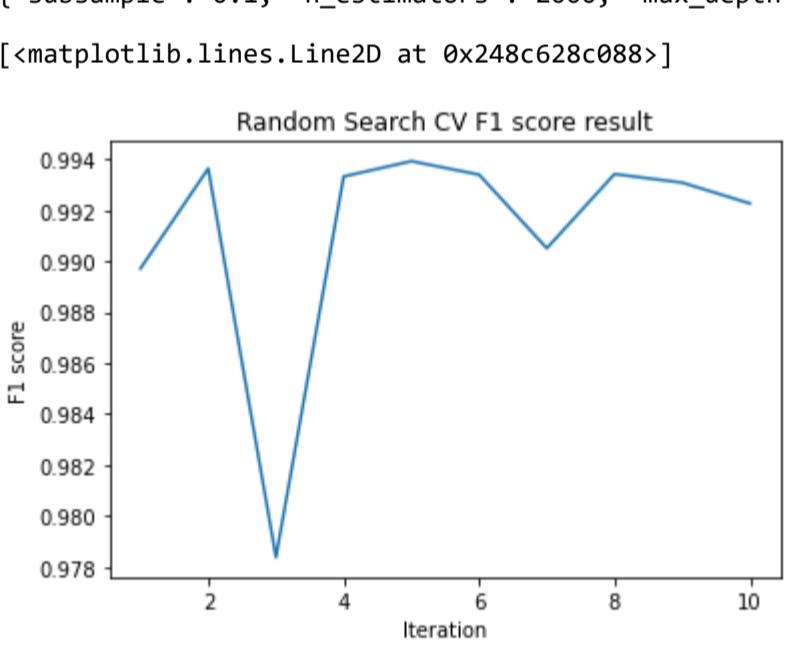
XG Boost Classifier

Median Imputed Data

```
In [99]: # XG Boost Classifier with Random Search CV for hyper parameter tuning
x_cfl=XGBClassifier()
params={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100, 500, 1000, 2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=params,n_jobs=-1, cv=3)
random_cfl.fit(X_train_median_os, y_train_median_os.values.ravel())
Out[99]: RandomizedSearchCV(cv=3,
                           estimator=XGBClassifier(base_score=None, booster=None,
                                                   colsample_bylevel=None,
                                                   colsample_bylevel=None,
                                                   colsample_bynode=None,
                                                   gamma=0, importance_type='gain',
                                                   gpu_id=None, interaction_constraints=None,
                                                   learning_rate=None,
                                                   max_delta_step=None, max_depth=None,
                                                   min_child_weight=None, missing=nan,
                                                   monotone_constraints=None,
                                                   n_estimators=100...),
                           n_estimators=100...),
                           num_parallel_tree=None,
                           random_state=None, reg_alpha=None,
                           reg_lambda=None,
                           scale_pos_weight=None,
                           subsample=None, tree_method=None,
                           validate_parameters=None,
                           verbosity=None),
                           n_jobs=1,
                           param_distributions={'colsample_bytree': [0.1, 0.3, 0.5, 1],
                                                   'learning_rate': [0.01, 0.03, 0.05, 0.1,
                                                       0.15, 0.2],
                                                   'max_depth': [3, 5, 10],
                                                   'n_estimators': [100, 500, 1000, 2000],
                                                   'subsample': [0.1, 0.3, 0.5, 1]})
```

```
In [100]: # Displaying mean test score which is the mean F1 score for each hyper parameter combination
labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)
```



```
In [101]: random_cfl.best_estimator_
Out[101]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                       colsample_bynode=1, colsample_bytree=0.1, gamma=0, gpu_id=-1,
                       importance_type='gain', interaction_constraints='',
                       learning_rate=0.01, max_delta_step=0, max_depth=10,
                       min_child_weight=1, missing=nan, monotone_constraints='()',
                       n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
                       reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                       tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [298]: from sklearn.calibration import CalibratedClassifierCV
x_cfl = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bynode=1, colsample_bytree=0.5, gamma=0, gpu_id=-1,
                      importance_type='gain', interaction_constraints='',
                      learning_rate=0.05, max_delta_step=0, max_depth=5,
                      min_child_weight=1, missing=nan, monotone_constraints='()',
                      n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
                      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                      tree_method='exact', validate_parameters=1, verbosity=None)
x_cfl.fit(X_train_median_os, y_train_median_os.values.ravel())
Out[298]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                       colsample_bynode=1, colsample_bytree=0.5, gamma=0, gpu_id=-1,
                       importance_type='gain', interaction_constraints='',
                       learning_rate=0.05, max_delta_step=0, max_depth=5,
                       min_child_weight=1, missing=nan, monotone_constraints='()',
                       n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
                       reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                       tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [299]: predict_y = x_cfl.predict(X_train_median_os)
print("Train F1 score:", f1_score(y_train_median_os, predict_y, average='macro'))
predict_y = x_cfl.predict(X_cv_median)
print("CV F1 score:", f1_score(y_cv, predict_y, average='macro'))
predict_y = x_cfl.predict(X_test_median)
print("Test F1 score:", f1_score(y_test, predict_y, average='macro'))
```

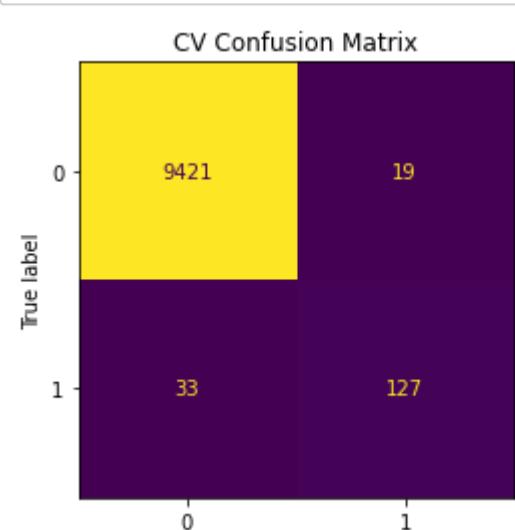
Train F1 score: 1.0

CV F1 score: 0.916565815063191

Test F1 score: 0.8781065869321554

With Calibrate classifier Train F1 score: 0.99 CV F1 score: 0.906 Test F1 score: 0.883

```
In [300]: plot1 = plot_confusion_matrix(x_cfl, X_cv_median, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')
plot2 = plot_confusion_matrix(x_cfl, X_test_median, y_test)
plot2.ax_.set_title('Test Confusion Matrix')
plt.show()
```



KNN Imputed Data

```
In [301]: # XG Boost Classifier with Random Search CV for hyper parameter tuning
x_cfl=XGBClassifier()

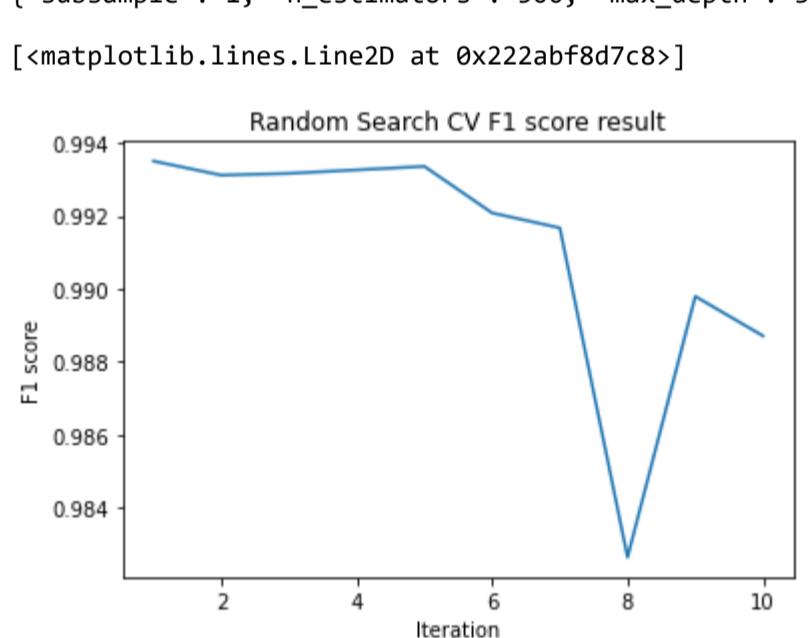
params={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100, 500, 1000, 2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl, param_distributions=params, n_jobs=-1, cv=3)
random_cfl.fit(X_train_knn_os, y_train_knn_os.values.ravel())
```

```
Out[301]: RandomizedSearchCV(cv=3,
                           estimator=XGBClassifier(base_score=None, booster=None,
                           colsample_bytree=None,
                           colsample_bylevel=None,
                           colsample_bynode=None,
                           colsample_bynode=None, gamma=None,
                           gpu_id=None, importance_type='gain',
                           interaction_constraints=None,
                           learning_rate=None,
                           max_delta_step=None, max_depth=None,
                           min_child_weight=None, missing=None,
                           monotone_constraints=None,
                           n_estimators=100,
                           n_jobs=-1,
                           num_parallel_trees=None,
                           random_state=None, reg_alpha=None,
                           reg_lambda=None,
                           scale_pos_weight=None,
                           subsample=None, tree_method=None,
                           validate_parameters=None,
                           verbosity=None),
                           n_jobs=-1,
                           param_distributions={'colsample_bytree': [0.1, 0.3, 0.5, 1],
                           'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2],
                           'max_depth': [3, 5, 10],
                           'n_estimators': [100, 500, 1000, 2000],
                           'subsample': [0.1, 0.3, 0.5, 1]})
```

```
In [302]: # Displaying mean test score which is the mean F1 score for each hyper parameter combination
```

```
labels = random_cfl.cv_results_['params']
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)
```



```
In [304]: random_cfl.best_estimator_
```

```
Out[304]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.3, gamma=0, gpu_id=-1,
importance_type='gain', interaction_constraints='',
learning_rate=0.2, max_delta_step=0, max_depth=10,
min_child_weight=1, missing=None, monotone_constraints=''),
n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.5,
tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [305]: x_cfl = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.3, gamma=0, gpu_id=-1,
importance_type='gain', interaction_constraints='',
learning_rate=0.2, max_delta_step=0, max_depth=10,
min_child_weight=1, missing=None, monotone_constraints=''),
n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.5,
tree_method='exact', validate_parameters=1, verbosity=None)
x_cfl.fit(X_train_knn_os,y_train_knn_os.values.ravel())
```

```
Out[305]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.3, gamma=0, gpu_id=-1,
importance_type='gain', interaction_constraints='',
learning_rate=0.2, max_delta_step=0, max_depth=10,
min_child_weight=1, missing=None, monotone_constraints=''),
n_estimators=2000, n_jobs=0, num_parallel_tree=1, random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.5,
tree_method='exact', validate_parameters=1, verbosity=None)
```

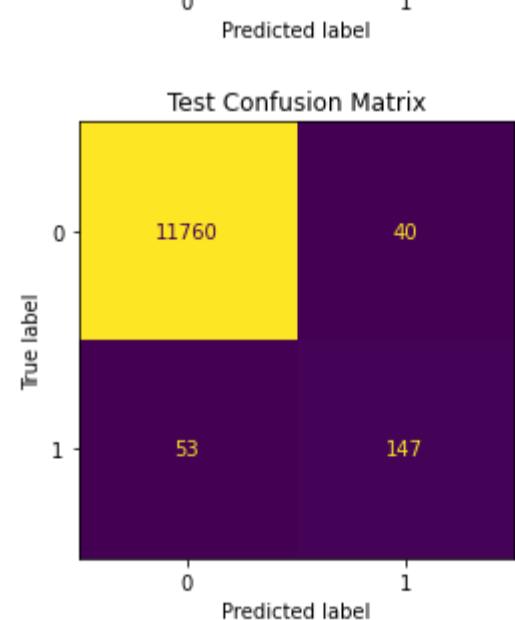
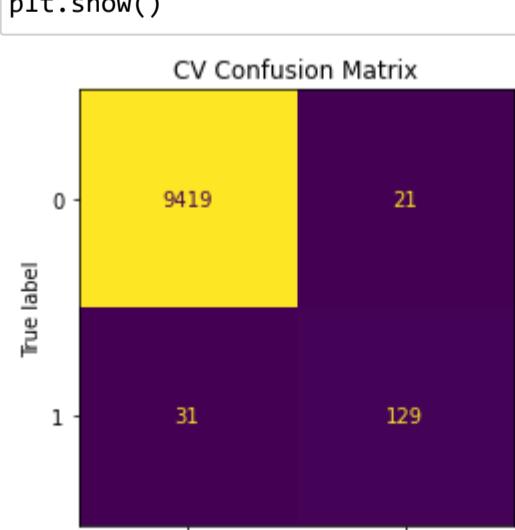
```
In [306]: predict_y = x_cfl.predict(X_train_knn_os)
print("Train F1 score: f1_score(y_train_knn_os, predict_y, average='macro')")
predict_y = x_cfl.predict(X_cv_knn_imputation)
print("CV F1 score: f1_score(y_cv, predict_y, average='macro')")
predict_y = x_cfl.predict(X_test_knn_imputation)
print("Test F1 score: f1_score(y_test, predict_y, average='macro')")
```

```
Train F1 score: 1.0
```

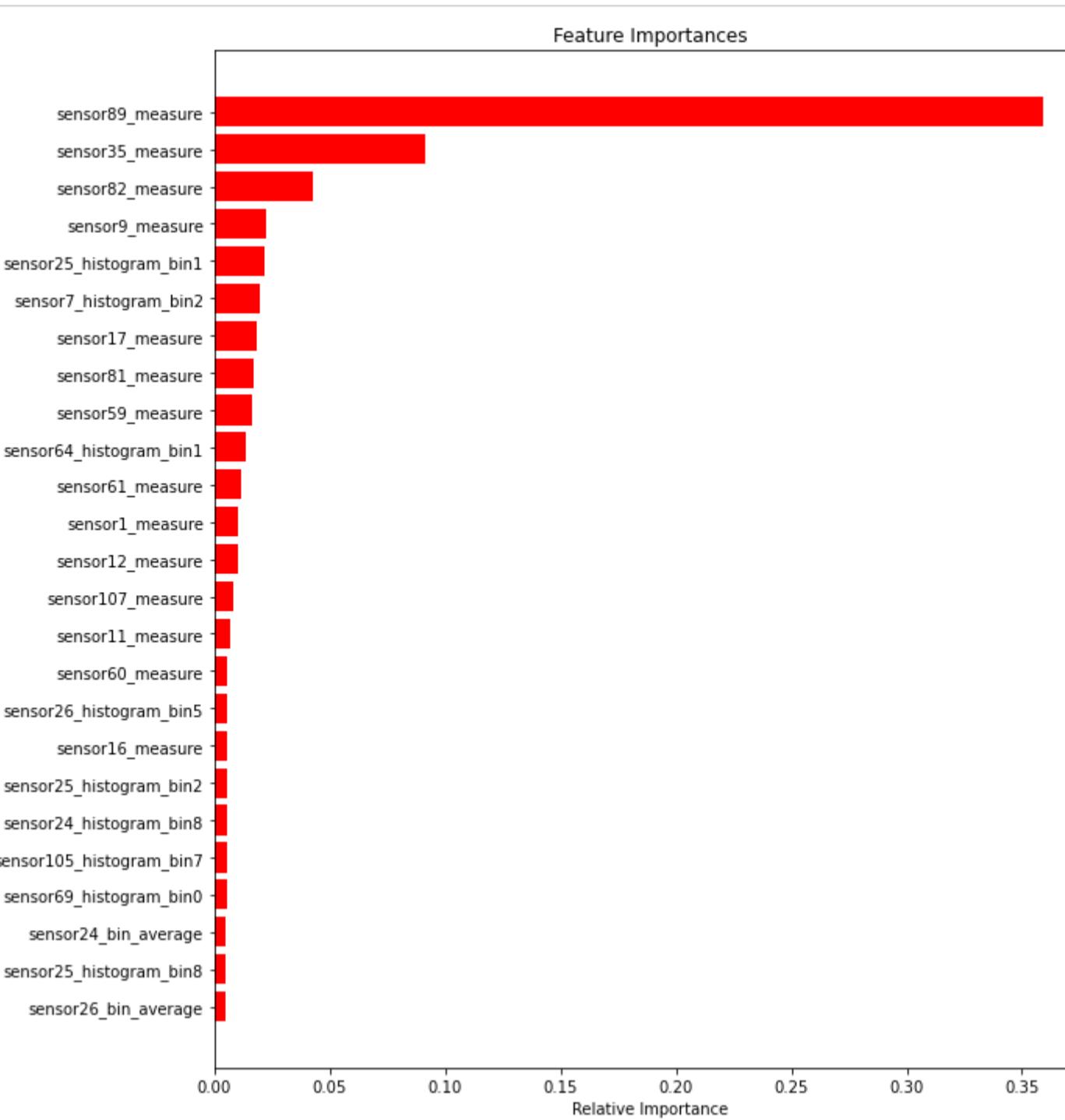
```
CV F1 score: 0.914752642339248
```

```
Test F1 score: 0.8778757070159422
```

```
In [307]: plot1 = plot_confusion_matrix(x_cfl, X_cv_knn_imputation, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')
plot2 = plot_confusion_matrix(x_cfl, X_test_knn_imputation, y_test)
plot2.ax_.set_title('Test Confusion Matrix')
plt.show()
```



```
In [308]: # Fetching the feature importances from the classifier and plotting the top 25 feature importances
features = X_train_knn.os.columns
importances = x_cfl.feature_importances_
indices = np.argsort(importances)[-25:]
plt.figure(figsize=(12, 2))
plt.title('Feature Importances')
plt.bar(range(len(indices)), importances[indices], color='r', align='center')
plt.yticks(range(len(indices)), [features[i] for i in indices])
plt.xlabel('Relative Importance')
plt.show()
```



Comments for XG Boost

- XG Boost Model provided slightly better results over the RF classifier mainly for the Majority class, whereas the classification of the minority class remains the same.
- We see that FE of Averaging Bin values has good feature importance and in the top 25 features of the 150+ features.
- We can try improving using stacking and voting classifier to see if they can classify better.

Since XG Boost provided good results, we can try for the non oversampled dataset with XG Boost for the knn based imputation to see if dataset without oversampling can do better than oversampling.

```
In [309]: mmx_scaler = MinMaxScaler()
columns = X_train_knn_imputation.columns
```

```
X_train_knn_wo_os = pd.DataFrame(mmx_scaler.fit_transform(X_train_knn_imputation))
X_cv_knn_wo_os = pd.DataFrame(mmx_scaler.transform(X_cv_knn_imputation))
X_test_knn_wo_os = pd.DataFrame(mmx_scaler.transform(X_test_knn_imputation))
```

```
X_train_knn_wo_os.columns = columns
X_cv_knn_wo_os.columns = columns
X_test_knn_wo_os.columns = columns
```

```
In [310]: x_cfl.XGBClassifier()
```

```
params = {
    'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2],
    'n_estimators': [100, 300, 1000, 2000],
    'max_depth': [3, 5, 10],
    'colsample_bytree': [0.1, 0.3, 0.5, 1],
    'subsample': [0.1, 0.3, 0.5, 1]
}
random_cfl.RandomizedSearchCV(x_cfl, params, n_jobs=-1, cv=3)
random_cfl.fit(X_train_knn_wo_os, y_train.values.ravel())
```

```
Out[310]: RandomizedSearchCV(cv=3,
```

```
        estimator=XGBClassifier(base_score=None, booster=None,
                                colsample_bylevel=None,
                                colsample_bynode=None,
                                colsample_bytree=None, gamma=None,
                                gpu_id=None, importance_type='gain',
                                interaction_constraints=None,
                                learning_rate=None,
                                max_delta_step=None, max_depth=None,
                                min_child_weight=None, missing='NaN',
                                monotone_constraints=None,
                                n_estimators=100...,
                                n_min_child_weight=None,
                                neta_lambda=None, reg_alpha=None,
                                random_state=None, reg_lambda=None,
                                scale_pos_weight=None,
                                subsample=None, tree_method=None,
                                validate_parameters=None,
                                verbosity=None),
        n_jobs=1,
        param_distributions={'colsample_bytree': [0.1, 0.3, 0.5, 1],
                             'learning_rate': [0.01, 0.03, 0.05, 0.1,
                                              0.15, 0.2],
                             'max_depth': [3, 5, 10],
                             'n_estimators': [100, 300, 1000, 2000],
                             'subsample': [0.1, 0.3, 0.5, 1]})
```

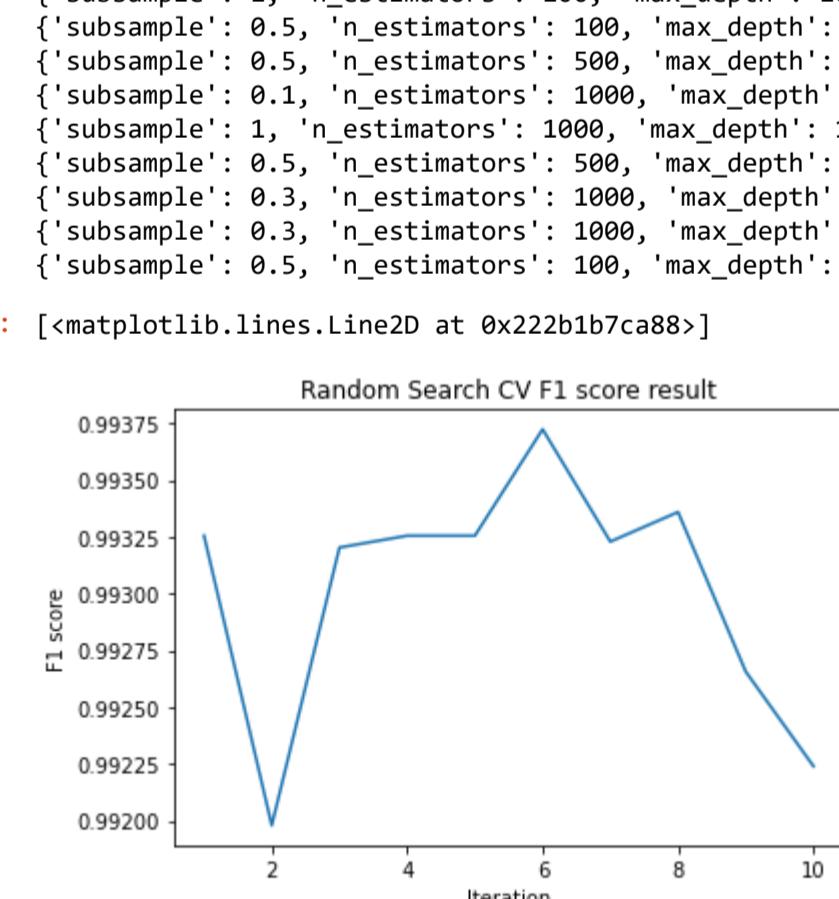
```
In [311]: labels = random_cfl.cv_results_['params']
```

```
x_axis = range(1,11)
y_axis = random_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])
```

```
plt.xlabel('Iteration')
```

```
plt.ylabel('F1 score')
```

```
plt.title('Random Search CV F1 score result')
```



```
In [312]: random_cfl.best_estimator_
```

```
Out[312]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                        colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=1,
                        importance_type='gain', interaction_constraints='',
                        learning_rate=0.1, max_delta_step=0, max_depth=10,
                        min_child_weight=1, missing='NaN', monotone_constraints='',
                        n_estimators=1000, n_jobs=0, num_parallel_tree=1, random_state=0,
                        reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                        tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [313]: x_cfl = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                            colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=1,
                            importance_type='gain', interaction_constraints='',
                            learning_rate=0.1, max_delta_step=0, max_depth=10,
                            min_child_weight=1, missing='NaN', monotone_constraints='')
                            n_estimators=1000, n_jobs=0, num_parallel_tree=1, random_state=0,
                            reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                            tree_method='exact', validate_parameters=1, verbosity=None)
```

```
Out[313]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                        colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=1,
                        importance_type='gain', interaction_constraints='',
                        learning_rate=0.1, max_delta_step=0, max_depth=10,
                        min_child_weight=1, missing='NaN', monotone_constraints=''),
                        n_estimators=1000, n_jobs=0, num_parallel_tree=1, random_state=0,
                        reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                        tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [314]: predict_y = x_cfl.predict(X_train_knn_wo_os)
print("Train F1 score: ", f1_score(y_train, predict_y, average='macro'))
predict_y = x_cfl.predict(X_cv_knn_wo_os)
print("CV F1 score: ", f1_score(y_cv, predict_y, average='macro'))
predict_y = x_cfl.predict(X_test_knn_wo_os)
print("Test F1 score: ", f1_score(y_test, predict_y, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.49579831932773105
Test F1 score: 0.49579831932773105
```

The non oversampled data provided did not give good results and the minority class has been misclassified significantly.

Next we use the Stacking and Voting classifiers in which we use the existing RandomSearch CV best results obtained from the above 3 models of KNN, RF and XG Boost.

Stacking Classifier

Median Imputed Data

```
In [315]: # Using 3 Stacking classifiers and one meta classifier with best hyperparameters obtained from the above
# individual classifiers hyperparameters tuning.

clf1 = RandomForestClassifier(max_depth=100, n_estimators=1000)
clf2 = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
                     colsample_bynode=1, colsample_bylevel=1,
                     importance_type='gain', interaction_constraints='',
                     learning_rate=0.05, max_delta_step=0, max_depth=10,
                     min_child_weight=1, monotone_constraints='()',
                     n_estimators=2000, n_jobs=-1, num_parallel_tree=1, random_state=0,
                     reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                     tree_method='exact', validate_parameters=1, verbosity=None)
clf3 = KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')
knn = KNeighborsClassifier(n_neighbors=5, weights='distance')

scf1 = StackingClassifier(classifiers=[clf1, clf2, clf3],
                         use_probas=False,
                         average_probas=False,
                         meta_classifier=knn)

scf1.fit(X_train_median_os, y_train_median_os.values.ravel())

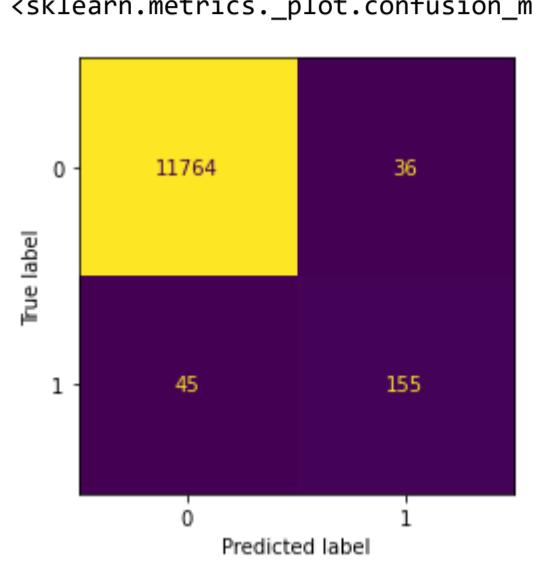
y_train_pred = scf1.predict(X_train_median_os)
print("Train F1 score:", f1_score(y_train_median_os, y_train_pred, average='macro'))

y_cv_pred = scf1.predict(X_cv_median)
print("CV F1 score:", f1_score(y_cv, y_cv_pred, average='macro'))

y_test_pred = scf1.predict(X_test_median)
print("Test F1 score:", f1_score(y_test, y_test_pred, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.9152901313082965
Test F1 score: 0.8947039898413183
```

```
In [316]: cm = confusion_matrix(y_test, y_test_pred, labels=[0,1])
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0,1])
disp.plot()
```



KNN imputed Data

```
In [318]: # Using 3 Stacking classifiers and one meta classifier with best hyperparameters obtained from the above
# individual classifiers hyperparameters tuning.

clf1 = RandomForestClassifier(max_depth=100, n_estimators=1000)
clf2 = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
                     colsample_bynode=1, colsample_bylevel=1,
                     importance_type='gain', interaction_constraints='',
                     learning_rate=0.1, max_delta_step=0, max_depth=10,
                     min_child_weight=1, monotone_constraints='()',
                     n_estimators=2000, n_jobs=-1, num_parallel_tree=1, random_state=0,
                     reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.5,
                     tree_method='exact', validate_parameters=1, verbosity=None)
clf3 = KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')
knn = KNeighborsClassifier(n_neighbors=5, weights='distance')

scf1 = StackingClassifier(classifiers=[clf1, clf2, clf3],
                         use_probas=False,
                         average_probas=False,
                         meta_classifier=knn)

scf1.fit(X_train_knn_os, y_train_knn_os.values.ravel())

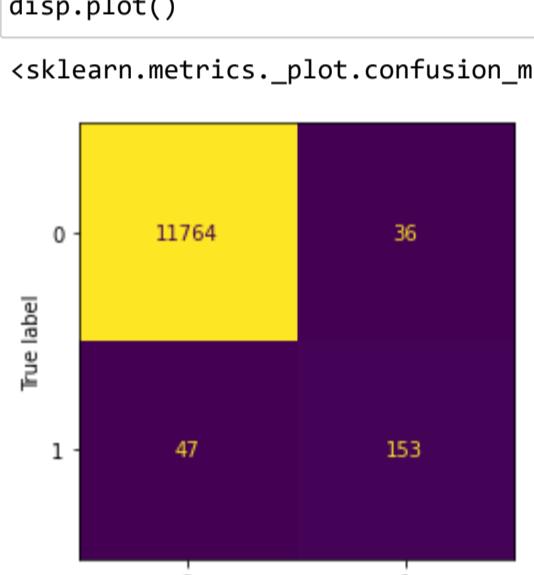
y_train_pred = scf1.predict(X_train_knn_os)
print("Train F1 score:", f1_score(y_train_knn_os, y_train_pred, average='macro'))

y_cv_pred = scf1.predict(X_cv_knn_imputation)
print("CV F1 score:", f1_score(y_cv,y_cv_pred, average='macro'))

y_test_pred = scf1.predict(X_test_knn_imputation)
print("Test F1 score:", f1_score(y_test, y_test_pred, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.9068513139633319
Test F1 score: 0.8915854042608
```

```
In [319]: cm = confusion_matrix(y_test, y_test_pred, labels=[0,1])
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0,1])
disp.plot()
```



The stacking classifier of the Median and KNN based imputation improved its performance over the XG Boost and provided better classifications for the majority class and thus there is a slight increase in the F1 score.

The Median based imputation performed marginally better over the KNN imputed dataset.

Voting Classifier

Median Imputed Data

```
In [102]: # Using 3 Voting classifiers and hard voting with best hyperparameters obtained from the above
# individual classifiers hyperparameters tuning.

rf = RandomForestClassifier(max_depth=100, n_estimators=2000)
xgb = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
                     colsample_bynode=1, colsample_bylevel=1,
                     importance_type='gain', interaction_constraints='',
                     learning_rate=0.1, max_delta_step=0, max_depth=10,
                     min_child_weight=1, monotone_constraints='()',
                     n_estimators=1000, n_jobs=-1, num_parallel_tree=1, random_state=0,
                     reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                     tree_method='exact', validate_parameters=1, verbosity=None)
knn = KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')

clf = VotingClassifier(estimators=[('rf', rf), ('xgb', xgb), ('knn', knn)], voting='hard', n_jobs=-1)

clf.fit(X_train_median_os, y_train_median_os.values.ravel())
```

```
Out[102]: VotingClassifier(estimators=[('rf',
                                         RandomForestClassifier(max_depth=100,
                                                               n_estimators=2000),
                                         ('xgb',
                                          XGBClassifier(base_score=0.5, booster='gbtree',
                                                        colsample_bytree=1,
                                                        colsample_bynode=1,
                                                        colsample_bylevel=1,
                                                        importance_type='gain',
                                                        interaction_constraints='',
                                                        learning_rate=0.1, max_delta_step=0,
                                                        max_depth=10, min_child_weight=1,
                                                        missing='nan',
                                                        monotone_constraints='()',
                                                        n_estimators=1000, n_jobs=2,
                                                        num_parallel_threads=1, random_state=0,
                                                        reg_alpha=0, reg_lambda=1,
                                                        scale_pos_weight=1, subsample=1,
                                                        tree_method='exact',
                                                        validate_parameters=1,
                                                        verbosity=None),
                                         ('knn',
                                          KNeighborsClassifier(metric='manhattan',
                                                               n_neighbors=4,
                                                               weights='distance'))),
                                         n_jobs=-1)
```

```
In [103]: y_train_pred = clf.predict(X_train_median_os)
print("Train F1 score:", f1_score(y_train_median_os, y_train_pred, average='macro'))

y_cv_pred = clf.predict(X_cv_median)
print("CV F1 score:", f1_score(y_cv,y_cv_pred, average='macro'))

y_test_pred = clf.predict(X_test_median)
print("Test F1 score:", f1_score(y_test, y_test_pred, average='macro'))
```

```
Train F1 score: 1.0
CV F1 score: 0.879966780135406
Test F1 score: 0.8696771250773793
```

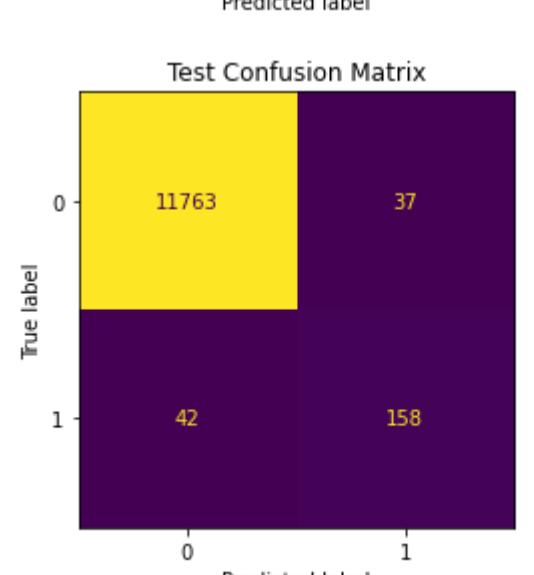
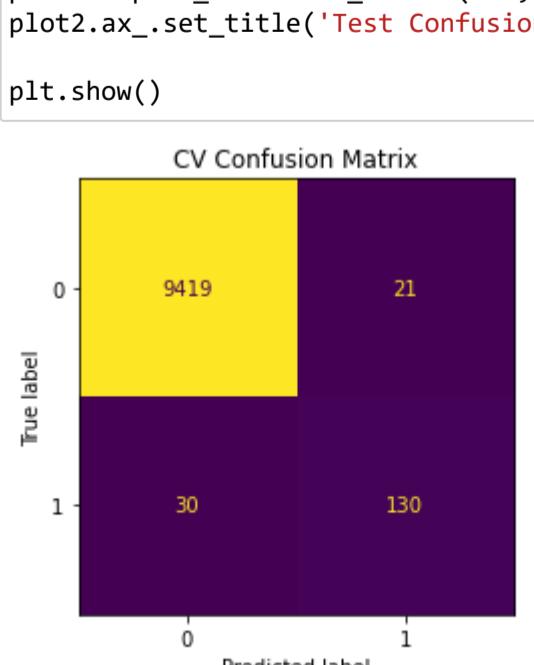
```
In [92]: with open('rf.pkl', 'wb') as file:
    pickle.dump(rf, file)

with open('xgb.pkl', 'wb') as file:
    pickle.dump(xgb, file)

with open('knn.pkl', 'wb') as file:
    pickle.dump(knn, file)

with open('clf_voting_classifier.pkl', 'wb') as file:
    pickle.dump(clf, file)
```

```
In [322]: plot1 = plot_confusion_matrix(clf, X_cv_median, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')
plot2 = plot_confusion_matrix(clf, X_test_median, y_test)
plot2.ax_.set_title('Test Confusion Matrix')
```



KNN imputed Data

```
In [323]: # Using 3 Voting classifiers and hard voting with best hyperparameters obtained from the above
# individual classifiers hyperparameters tuning.

rf = RandomForestClassifier(max_depth=100, n_estimators=1000)
xgb = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
                     colsample_bynode=1, colsample_bylevel=1,
                     gamma=0, gpu_id=-1,
                     importance_type='gain', interaction_constraints='',
                     learning_rate=0.2, max_delta_step=0,
                     min_child_weight=1, monotone_constraints='',
                     n_estimators=2000, n_jobs=-1, num_parallel_tree=1, random_state=0,
                     reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.5,
                     tree_method='exact', validate_parameters=1, verbosity=None)
knn = KNeighborsClassifier(metric='manhattan', n_neighbors=4, weights='distance')

clf = VotingClassifier(estimators=[('rf', rf), ('xgb', xgb), ('knn', knn)], voting='hard', n_jobs=-1)

clf.fit(X_train_knn_os, y_train_knn_os.values.ravel())

Out[323]: VotingClassifier(estimators=[('rf',
                                         RandomForestClassifier(max_depth=100,
                                                               n_estimators=1000)),
                                         ('xgb',
                                         XGBClassifier(base_score=0.5, booster='gbtree',
                                                       colsample_bytree=1,
                                                       colsample_bynode=1,
                                                       colsample_bylevel=1,
                                                       gamma=0,
                                                       gpu_id=-1, importance_type='gain',
                                                       interaction_constraints='',
                                                       learning_rate=0.2, max_delta_step=0,
                                                       min_child_weight=1, monotone_constraints='',
                                                       n_estimators=2000, n_jobs=-1,
                                                       num_parallel_tree=1, random_state=0,
                                                       reg_alpha=0, reg_lambda=1,
                                                       scale_pos_weight=1, subsample=0.5,
                                                       tree_method='exact',
                                                       validate_parameters=1,
                                                       verbosity=None),
                                         ('knn',
                                         KNeighborsClassifier(metric='manhattan',
                                                               n_neighbors=4,
                                                               weights='distance'))],
                                         n_jobs=-1)
```

```
In [324]: y_train_pred = clf.predict(X_train_knn_os)
print("Train F1 score:", f1_score(y_train_knn_os, y_train_pred, average="macro"))

y_cv_pred = clf.predict(X_cv_knn_imputation)
print("CV F1 score:", f1_score(y_cv_knn_imputation, y_cv_pred, average="macro"))

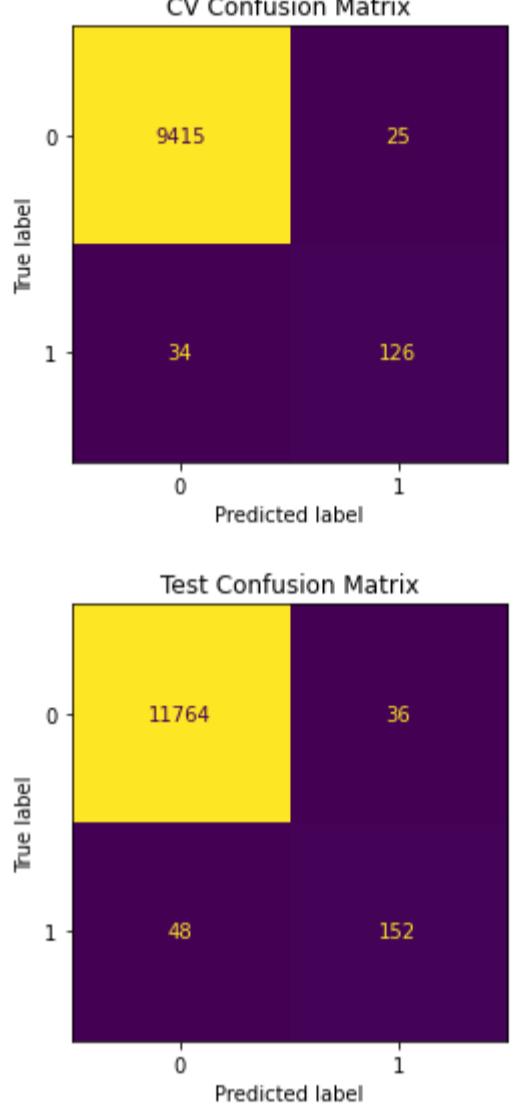
y_test_pred = clf.predict(X_test_knn_imputation)
print("Test F1 score:", f1_score(y_test, y_test_pred, average="macro"))

Train F1 score: 1.0
CV F1 score: 0.9835829398146769
Test F1 score: 0.8995382055129
```

```
In [325]: plot1 = plot_confusion_matrix(clf, X_cv_knn_imputation, y_cv)
plot1.ax_.set_title('CV Confusion Matrix')

plot2 = plot_confusion_matrix(clf, X_test_knn_imputation, y_test)
plot2.ax_.set_title('Test Confusion Matrix')

plot.show()
```



Voting Classifiers also provided similar results as the stacking classifier. The Median based imputation + Voting classifier provided the best result till now with an test F1 score of 0.898.

Ensemble Model

with Decision Tree and Logistic Regression as the metaclassifier

Without feature engineering data

```
In [322]: X_train, X_ensemble_test, y_train, y_ensemble_test = train_test_split(data_df, y_true, stratify=y_true, test_size=0.2)
dataset_d1, dataset_d2, y_d1, y_d2 = train_test_split(X_train, y_train, stratify=y_train, test_size=0.5)
print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 149) (24000, 149) (12000, 149)

In [333]: high_correlated_columns = ['sensor45_measure', 'sensor32_measure', 'sensor46_measure', 'sensor47_measure', 'sensor67_measure']
for column in high_correlated_columns:
    if column in X_train.columns:
        dataset_d1 = dataset_d1.drop([column].axis=1)
        dataset_d2 = dataset_d2.drop([column].axis=1)
        X_ensemble_test = X_ensemble_test.drop([column], axis=1)

dataset_d1 = dataset_d1.drop(['sensor04_histogram_bin5'], axis=1)
dataset_d2 = dataset_d2.drop(['sensor04_histogram_bin5'], axis=1)
X_ensemble_test = X_ensemble_test.drop(['sensor04_histogram_bin5'], axis=1)

In [335]: knn_imputer = KNNImputer(weights="distance")
columns = dataset_d1.columns
dataset_d1 = knn_imputer.fit_transform(dataset_d1)
dataset_d2 = knn_imputer.transform(dataset_d2)
X_ensemble_test = knn_imputer.transform(X_ensemble_test)

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 143) (24000, 143) (12000, 143)

In [336]: dataset_d1 = pd.DataFrame(dataset_d1)
dataset_d2 = pd.DataFrame(dataset_d2)
X_ensemble_test = pd.DataFrame(X_ensemble_test)

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns

In [337]: dataset_d1
```

	sensor1_measure	sensor3_measure	sensor4_measure	sensor5_measure	sensor6_measure	sensor7_histogram_bin2	sensor7_histogram_bin3	sensor7_histogram_bin4	sensor7_histogram_bin5	sensor7_histogram_bin6	...	sensor105_histogram_bin1	sensor105_histogram_bin2	sensor105_histogram_bin3	sensor105_histogram_bin4	sensor105_histogram_bin5	sensor105_histogram_bin6	sensor105_histogram_bin7	sensor105_hist
0	39222.0	2.130706e+09	624.000000	0.0	0.0	0.0	47830.0	554420.0	1195392.0	...	425902.0	320650.0	160140.0	299866.0	263108.0	222682.0	201290.0	201290.0	
1	1590.0	3.400000e+01	34.000000	0.0	0.0	0.0	47802.0	24822.0	52180.0	4924.0	12422.0	5054.0	2184.0	2820.0	3170.0	8006.0	44652.0		
2	1312.0	0.000000e+00	0.000000	0.0	0.0	0.0	0.0	1010.0	71870.0	3838.0	9654.0	2256.0	622.0	2364.0	2378.0	1848.0	43536.0		
3	61336.0	2.580000e+02	226.000000	0.0	0.0	0.0	179042.0	150222.0	2306598.0	...	956472.0	817226.0	506226.0	910222.0	609818.0	89806.0	30416.0		
4	918058.0	6.536000e+03	0.000000	0.0	0.0	0.0	728572.0	11616492.0	16716256.0	6624548.0	...	3514750.0	2730552.0	1572668.0	2602968.0	2904338.0	1928112.0	1789334.0	
...		
23995	596132.0	0.000000e+00	2748.390879	0.0	0.0	0.0	384844.0	1192200.0	2169028.0	...	7719156.0	5437924.0	2867092.0	5779214.0	4861518.0	3532022.0	1900662.0		
23996	10776.0	4.740000e+02	87.035008	0.0	0.0	0.0	21896.0	161356.0	322290.0	...	53134.0	45374.0	30380.0	318468.0	70222.0	29978.0	11314.0		
23997	1348.0	4.800000e+01	44.000000	0.0	0.0	0.0	70.0	11344.0	70956.0	...	14462.0	9372.0	7162.0	39224.0	9470.0	166.0	10.0		
23998	38550.0	2.130706e+09	756.000000	0.0	0.0	0.0	178.0	36977.0	1697408.0	415174.0	216968.0	116980.0	248000.0	214792.0	172696.0	13666.0			
23999	60150.0	6.060000e+02	470.000000	0.0	0.0	0.0	1494.0	40074.0	2768462.0	...	739992.0	491072.0	212864.0	489998.0	747832.0	548846.0	204828.0		

24000 rows x 143 columns

```
In [338]: mmnx_scaler = MinMaxScaler()
columns = dataset_d1.columns

dataset_d1 = pd.DataFrame(mmnx_scaler.fit_transform(dataset_d1))
dataset_d2 = pd.DataFrame(mmnx_scaler.transform(dataset_d2))
X_ensemble_test = pd.DataFrame(mmnx_scaler.transform(X_ensemble_test))

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [337]: dataset_d1
```

	sensor1_measure	sensor3_measure	sensor4_measure	sensor5_measure	sensor6_measure	sensor7_histogram_bin2	sensor7_histogram_bin3	sensor7_histogram_bin4	sensor7_histogram_bin5	sensor7_histogram_bin6	...	sensor105_histogram_bin1	sensor105_histogram_bin2	sensor105_histogram_bin3	sensor105_histogram_bin4	sensor105_histogram_bin5	sensor105_histogram_bin6	sensor105_histogram_bin7	sensor105_hist
0	39222.0	2.130706e+09	624.000000	0.0	0.0	0.0	47830.0	554420.0	1195392.0	...	425902.0	320650.0	160140.0	299866.0	263108.0	222682.0	201290.0	201290.0	
1	1590.0	3.400000e+01	34.000000	0.0	0.0	0.0	47802.0	24822.0	52180.0	4924.0	12422.0	5054.0	2184.0	2820.0	3170.0	8006.0	44652.0		
2	1312.0	0.000000e+00	0.000000	0.0	0.0	0.0	0.0	1010.0	71870.0	3838.0	9654.0	2256.0	622.0	2364.0	2378.0	1848.0	43536.0		
3	61336.0	2.580000e+02	226.000000	0.0	0.0	0.0	179042.0	150222.0	2306598.0	...	956472.0	817226.0	506226.0	910222.0	609818.0	89806.0	30416.0		
4	918058.0	6.536000e+03	0.000000	0.0	0.0	0.0	728572.0	11616492.0	16716256.0	6624548.0	...	3514750.0	2730552.0	1572668.0	2602968.0	2904338.0	1928112.0	1789334.0	
...		
23995	596132.0	0.000000e+00	2748.390879	0.0	0.0	0.0	384844.0	1192200.0	2169028.0	...	7719156.0	5437924.0	2867092.0	5779214.0	4861518.0	3532022.0	1900662.0		
23996	10776.0	4.740000e+02	87.035008	0.0	0.0	0.0	21896.0	161356.0	322290.0	...	53134.0	45374.0	30380.0	318468.0	70222.0	29978.0	11314.0		
23997	1348.0	4.800000e+01	44.000000	0.0	0.0	0.0	70.0	11344.0	70956.0	...	14462.0	9372.0	7162.0	39224.0	9470.0	166.0	10.0		
23998	38550.0	2.130706e+09	756.000000	0.0	0.0	0.0	178.0	36977.0	1697408.0	415174.0	216968.0	116980.0	248000.0	214792.0	172696.0	13666.0			
23999	60150.0	6.060000e+02	470.000000	0.0	0.0	0.0	1494.0	40074.0	2768462.0	...	739992.0	491072.0	212864.0	489998.0	747832.0	548846.0	204828.0		

24000 rows x 143 columns

```
In [338]: mmnx_scaler = MinMaxScaler()
columns = dataset_d1.columns

dataset_d1 = pd.DataFrame(mmnx_scaler.fit_transform(dataset_d1))
dataset_d2 = pd.DataFrame(mmnx_scaler.transform(dataset_d2))
X_ensemble_test = pd.DataFrame(mmnx_scaler.transform(X_ensemble_test))

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [337]: dataset_d1
```

	sensor1_measure	sensor3_measure	sensor4_measure	sensor5_measure	sensor6_measure	sensor7_histogram_bin2	sensor7_histogram_bin3	sensor7_histogram_bin4	sensor7_histogram_bin5	sensor7_histogram_bin6	...	sensor105_histogram_bin1	sensor105_histogram_bin2	sensor105_histogram_bin3	sensor105_histogram_bin4	sensor105_histogram_bin5	sensor105_histogram_bin6	sensor105_histogram_bin7	sensor105_hist
0	39222.0	2.130706e+09	624.000000	0.0	0.0	0.0	47830.0	554420.0	1195392.0	...	425902.0	320650.0	160140.0	299866.0	263108.0	222682.0	201290.0	201290.0	
1	1590.0	3.400000e+01	34.000000	0.0	0.0	0.0	47802.0	24822.0	52180.0	4924.0	12422.0	5054.0	2184.0	2820.0	3170.0	800			

```
In [355]: # Number of models here is a hyper parameter and has been tried with multiple values and the results are below
num_of_models = 75
model_list, df = ensemble_model_DT(num_of_models, dataset_d1, y_d1, dataset_d2, y_d2)

Generating Samples Done
Fitting Base Models Done
Predict on Dataset 2 Done
Formed Meta Dataset Done
```

```
In [357]: # Logistic Regression as the meta classifier and it is hyper parameterized with the GridSearch CV
x_cfl = LogisticRegression()

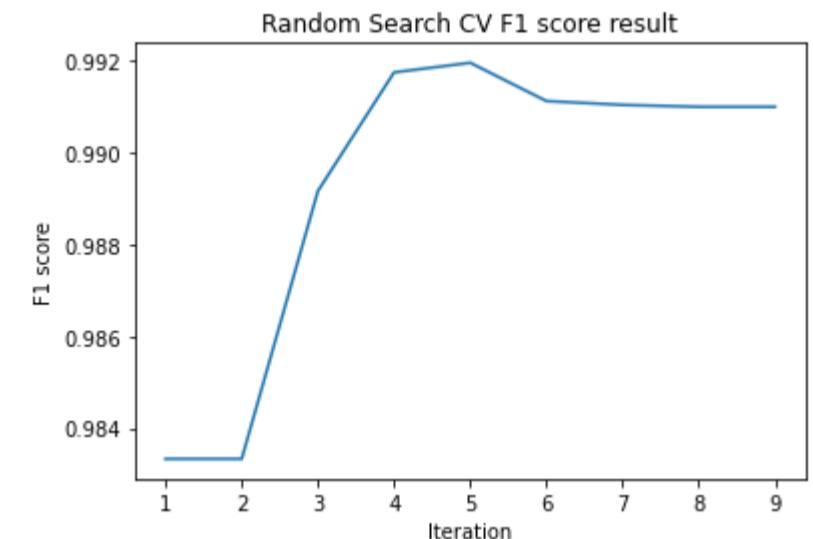
params={
    'C':[10 ** x for x in range(-5, 4)]
}
grid_cfl=GridSearchCV(x_cfl, param_grid=params, n_jobs=-1, scoring='f1')
grid_cfl.fit(df, y_d2.values.ravel())

Out[357]: GridSearchCV(estimator=LogisticRegression(), n_jobs=-1,
param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,
1000]})
```

```
In [358]: # plotting the various mean f1 scores
labels = grid_cfl.cv_results_['params']
x_axis = range(1,10)
y_axis = grid_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)

print("BEST MODEL -- ", grid_cfl.best_estimator_)
```



```
In [359]: x_cfl = LogisticRegression(C=0.1)
x_cfl.fit(df,y_d2.values.ravel())

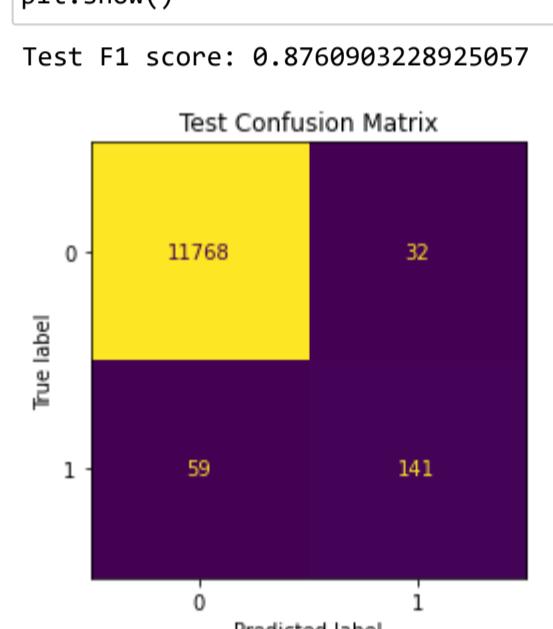
Out[359]: LogisticRegression(C=0.1)
```

```
In [360]: # Predicting the intermediate output for the test dataset before giving it to the meta classifier
def test_dataset(num_of_models, model_list, X_ensemble_test):
    test_output_list = []
    test_df = []

    # Predict output for the test set
    for i in range(num_of_models):
        model = model_list[i]
        output = model.predict(X_ensemble_test)
        test_output_list.append(output)

    # Zip the output of the test set two form a dataset
    for i in range(len(test_output_list)):
        output = []
        for j in range(num_of_models):
            output.append(test_output_list[j][i])
        test_df.append(output)
    test_df = pd.DataFrame(test_df)
    return test_df
```

```
In [361]: test_df = test_dataset(num_of_models, model_list, X_ensemble_test)
y_test_pred = x_cfl.predict(test_df)
print("Test F1 score:", f1_score(y_ensemble_test, y_test_pred, average='macro'))
```



Results

- 10 DT Models - 0.853355893965031
- 20 DT Models - 0.8513105921677959
- 50 DT Models - 0.8631384259778475
- 75 DT Models - 0.8760903228925057

Median Imputed Data without FE

```
In [362]: X_train, X_ensemble_test, y_train, y_ensemble_test = train_test_split(data_df, y_true, stratify=y_true, test_size=0.2)

dataset_d1, dataset_d2, y_d1, y_d2 = train_test_split(X_train, y_train, stratify=y_train, test_size=0.5)

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 149) (24000, 149) (12000, 149)
```

```
In [363]: high_correlated_columns = ['sensor45_measure', 'sensor32_measure', 'sensor46_measure', 'sensor47_measure', 'sensor67_measure']

for column in high_correlated_columns:
    if column in X_train.columns:
        dataset_d1 = dataset_d1.drop([column],axis=1)
        dataset_d2 = dataset_d2.drop([column],axis=1)
        X_ensemble_test = X_ensemble_test.drop([column],axis=1)

dataset_d1 = dataset_d1.drop(['sensor64_histogram_bin5'],axis=1)
dataset_d2 = dataset_d2.drop(['sensor64_histogram_bin5'],axis=1)
X_ensemble_test = X_ensemble_test.drop(['sensor64_histogram_bin5'],axis=1)
```

```
In [364]: median_values = dataset_d1.median()
columns = dataset_d1.columns
dataset_d1 = dataset_d1.fillna(median_values)
dataset_d2 = dataset_d2.fillna(median_values)
X_ensemble_test = X_ensemble_test.fillna(median_values)

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 143) (24000, 143) (12000, 143)
```

```
In [365]: dataset_d1 = pd.DataFrame(dataset_d1)
dataset_d2 = pd.DataFrame(dataset_d2)
X_ensemble_test = pd.DataFrame(X_ensemble_test)

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [366]: mmxx_scaler = MinMaxScaler()
columns = dataset_d1.columns

dataset_d1 = pd.DataFrame(mmxx_scaler.fit_transform(dataset_d1))
dataset_d2 = pd.DataFrame(mmxx_scaler.transform(dataset_d2))
X_ensemble_test = pd.DataFrame(mmxx_scaler.transform(X_ensemble_test))

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [383]: # Number of models here is a hyper parameter and has been tried with multiple values and the results are below
num_of_models = 75
model_list, df = ensemble_model_DT(num_of_models, dataset_d1, y_d1, dataset_d2, y_d2)

Generating Samples Done
Fitting Base Models Done
Predict on Dataset 2 Done
Formed Meta Dataset Done
```

```
In [384]: # Logistic Regression as the meta classifier and it is hyper parameterized with the GridSearch CV
x_cfl = LogisticRegression()

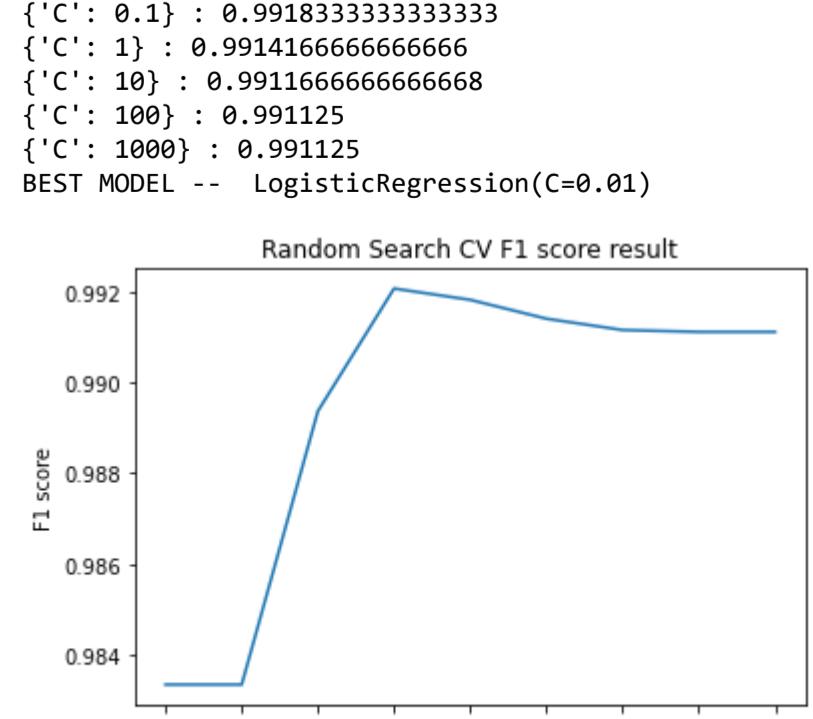
params={
    'C':[10 ** x for x in range(-5, 4)]
}
grid_cfl=GridSearchCV(x_cfl, param_grid=params, n_jobs=-1, scoring='f1')
grid_cfl.fit(df, y_d2.values.ravel())

Out[384]: GridSearchCV(estimator=LogisticRegression(), n_jobs=-1,
param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,
1000]})
```

```
In [385]: # plotting the various mean f1 scores
labels = grid_cfl.cv_results_['params']
x_axis = range(1,10)
y_axis = grid_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)

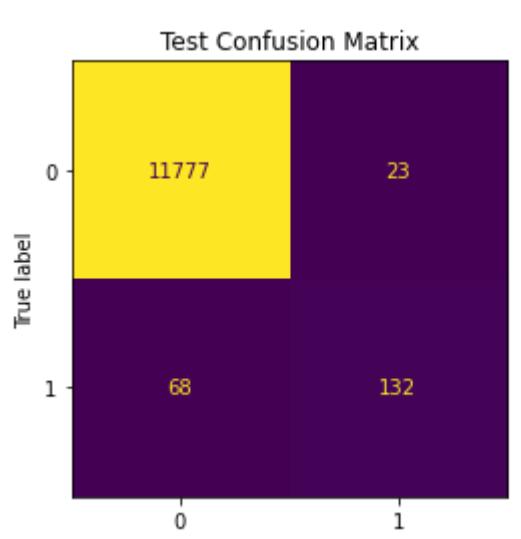
print("BEST MODEL -- ", grid_cfl.best_estimator_)
```



```
In [392]: x_cfl = LogisticRegression(C=0.1)
x_cfl.fit(df,y_d2.values.ravel())

Out[392]: LogisticRegression(C=0.1)
```

```
In [393]: test_df = test_dataset(num_of_models, model_list, X_ensemble_test)
y_test_pred = x_cfl.predict(test_df)
print("Test F1 score:", f1_score(y_ensemble_test, y_test_pred, average='macro'))
plot1 = plot_confusion_matrix(x_cfl, test_df, y_ensemble_test)
plot1.ax_.set_title('Test Confusion Matrix')
plot1.show()
Test F1 score: 0.8699066890239725
```



- 10 DT Models - 0.8399217868838921
- 20 DT Models - 0.889436052748646
- 50 DT Models - 0.8778061312810195
- 75 DT Models - 0.8699066890239725

Comments

From both the ensemble models we see that they are not very close to the Voting Classifiers when Feature Engineering has not been applied. This shows that our Feature Engineering has contributed to the final results and improving the F1 score.

Ensemble model

with feature engineered data for KNN imputed dataset

```
In [394]: X_train, X_ensemble_test, y_train, y_ensemble_test = train_test_split(data_df, y_true, stratify=y_true, test_size=0.2)
dataset_d1, dataset_d2, y_d1, y_d2 = train_test_split(X_train, y_train, stratify=y_train, test_size=0.5)
print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 149) (24000, 149) (24000, 149)

In [395]: high_correlated_columns = ['sensor45_measure', 'sensor32_measure', 'sensor46_measure', 'sensor47_measure', 'sensor67_measure']
for col in high_correlated_columns:
    if col in X_train.columns:
        dataset_d1 = dataset_d1.drop([column], axis=1)
        dataset_d2 = dataset_d2.drop([column], axis=1)
        X_ensemble_test = X_ensemble_test.drop([column], axis=1)

dataset_d1 = dataset_d1.drop(['sensor64_histogram_bin5'], axis=1)
dataset_d2 = dataset_d2.drop(['sensor64_histogram_bin5'], axis=1)
X_ensemble_test = X_ensemble_test.drop(['sensor64_histogram_bin5'], axis=1)
```

```
In [396]: knn_imputer = KNNImputer(weights="distance")
columns = dataset_d1.columns
dataset_d1 = knn_imputer.fit_transform(dataset_d1)
dataset_d2 = knn_imputer.transform(dataset_d2)
X_ensemble_test = knn_imputer.transform(X_ensemble_test)

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 143) (24000, 143) (24000, 143)
```

```
In [397]: dataset_d1 = pd.DataFrame(dataset_d1)
dataset_d2 = pd.DataFrame(dataset_d2)
X_ensemble_test = pd.DataFrame(X_ensemble_test)

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [398]: T_SVD = TruncatedSVD(n_components=4, n_iter=20)

dataset1_TSVD = T_SVD.fit_transform(dataset_d1)
dataset2_TSVD = T_SVD.transform(dataset_d2)
X_test_TSVD = T_SVD.fit_transform(X_ensemble_test)

print(dataset1_TSVD.shape, dataset2_TSVD.shape, X_test_TSVD.shape)
(24000, 4) (24000, 4) (24000, 4)
```

```
In [399]: for i in range(4):
    dataset_d1["TSVD"+str(i)] = dataset1_TSVD[:,i]
    dataset_d2["TSVD"+str(i)] = dataset2_TSVD[:,i]
    X_ensemble_test["TSVD"+str(i)] = X_test_TSVD[:,i]

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)
(24000, 147) (24000, 147) (24000, 147)
```

```
In [401]: column_names = dict()

for column in dataset_d1.columns:
    value = column.split('_')
    if value[1] == 'histogram':
        if value[0] not in column_names:
            column_names[value[0]] = []
        column_names[value[0]].append(column)

print(column_names)
```

```
{'sensor7': ['sensor7_histogram_bin2', 'sensor7_histogram_bin4', 'sensor7_histogram_bin5', 'sensor7_histogram_bin6', 'sensor7_histogram_bin8', 'sensor7_histogram_bin9'], 'sensor24': ['sensor24_histogram_bin5', 'sensor24_histogram_bin7', 'sensor24_histogram_bin9'], 'sensor25': ['sensor25_histogram_bin1', 'sensor25_histogram_bin3', 'sensor25_histogram_bin5', 'sensor25_histogram_bin7', 'sensor25_histogram_bin9'], 'sensor26': ['sensor26_histogram_bin1', 'sensor26_histogram_bin3', 'sensor26_histogram_bin5', 'sensor26_histogram_bin7', 'sensor26_histogram_bin9'], 'sensor27': ['sensor27_histogram_bin1', 'sensor27_histogram_bin3', 'sensor27_histogram_bin5', 'sensor27_histogram_bin7', 'sensor27_histogram_bin9'], 'sensor36': ['sensor36_histogram_bin1', 'sensor36_histogram_bin3', 'sensor36_histogram_bin5', 'sensor36_histogram_bin7', 'sensor36_histogram_bin9'], 'sensor64': ['sensor64_histogram_bin1', 'sensor64_histogram_bin2', 'sensor64_histogram_bin4', 'sensor64_histogram_bin6', 'sensor64_histogram_bin7', 'sensor64_histogram_bin8'], 'sensor65': ['sensor65_histogram_bin1', 'sensor65_histogram_bin3', 'sensor65_histogram_bin5', 'sensor65_histogram_bin7', 'sensor65_histogram_bin9'], 'sensor66': ['sensor66_histogram_bin1', 'sensor66_histogram_bin3', 'sensor66_histogram_bin5', 'sensor66_histogram_bin7', 'sensor66_histogram_bin9'], 'sensor67': ['sensor67_histogram_bin1', 'sensor67_histogram_bin3', 'sensor67_histogram_bin5', 'sensor67_histogram_bin7', 'sensor67_histogram_bin9'], 'sensor68': ['sensor68_histogram_bin1', 'sensor68_histogram_bin3', 'sensor68_histogram_bin5', 'sensor68_histogram_bin7', 'sensor68_histogram_bin9'], 'sensor69': ['sensor69_histogram_bin1', 'sensor69_histogram_bin3', 'sensor69_histogram_bin5', 'sensor69_histogram_bin7', 'sensor69_histogram_bin9'], 'sensor105': ['sensor105_histogram_bin1', 'sensor105_histogram_bin3', 'sensor105_histogram_bin5', 'sensor105_histogram_bin7', 'sensor105_histogram_bin9'], 'sensor106': ['sensor106_histogram_bin1', 'sensor106_histogram_bin3', 'sensor106_histogram_bin5', 'sensor106_histogram_bin7', 'sensor106_histogram_bin9']}

In [402]: dataset_d1 = add_average_bin_feature(dataset_d1, column_names)
dataset_d2 = add_average_bin_feature(dataset_d2, column_names)
X_ensemble_test = add_average_bin_feature(X_ensemble_test, column_names)
```

```
In [403]: mmnx_scaler = MinMaxScaler()
columns = dataset_d1.columns

dataset_d1 = pd.DataFrame(mmxn_scaler.fit_transform(dataset_d1))
dataset_d2 = pd.DataFrame(mmxn_scaler.transform(dataset_d2))
X_ensemble_test = pd.DataFrame(mmxn_scaler.transform(X_ensemble_test))

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns
```

```
In [448]: # Number of models here is a hyper parameter and has been tried with multiple values and the results are below
num_of_models = 75
model_list, df = ensemble_model_DT(num_of_models, dataset_d1, y_d1, dataset_d2, y_d2)

Generating Samples Done
Fitting Base Models Done
Predict on Dataset 2 Done
Formed Meta Dataset Done
```

```
In [442]: # Logistic Regression as the meta classifier and it is hyper parameterized with the GridSearch CV
```

```
x_cfl = LogisticRegression()
params = {
    'C': [10 ** x for x in range(-5, 4)]
}
grid_cfl = GridSearchCV(x_cfl, param_grid=params, n_jobs=-1, scoring='f1')
grid_cfl.fit(df, y_d2.values.ravel())
Out[442]: GridSearchCV(estimator=LogisticRegression(), param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000]})
```

```
In [443]: # plotting the various mean f1 scores
```

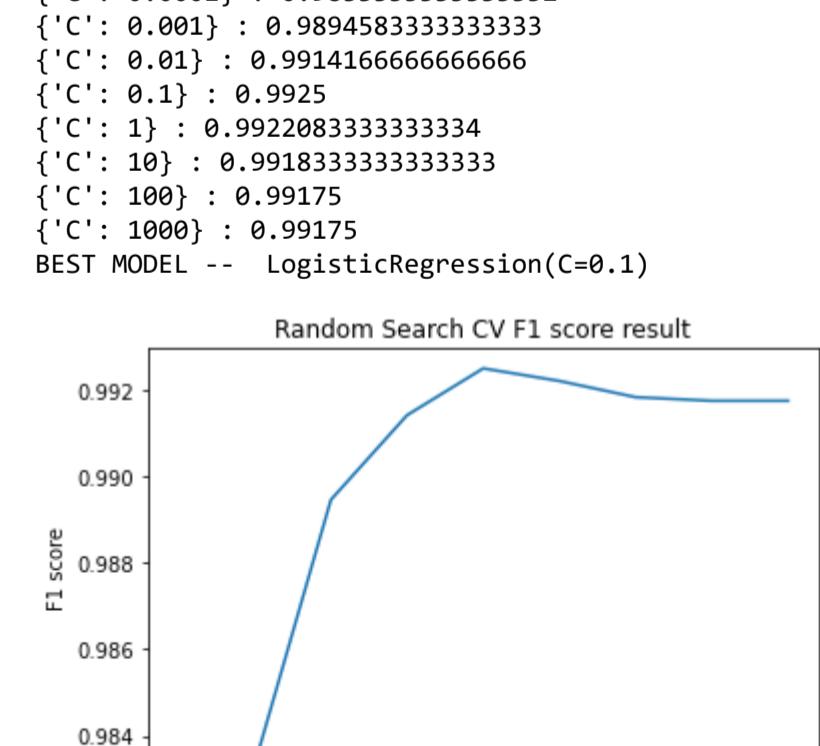
```
labels = grid_cfl.cv_results_['params']
x_axis = range(1,10)
y_axis = grid_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label, ":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)
```

```
print("BEST MODEL -- ", grid_cfl.best_estimator_)
```

```
{'C': 1e-05} : 0.9031353333333333
{'C': 0.0001} : 0.9033333333333333
{'C': 0.001} : 0.9045823232323232
{'C': 0.01} : 0.9914166666666666
{'C': 0.025} : 0.9925
{'C': 0.1} : 0.9922883333333333
{'C': 0.1001} : 0.9918333333333333
{'C': 1.0001} : 0.99175
BEST MODEL -- LogisticRegression(C=0.1)
```

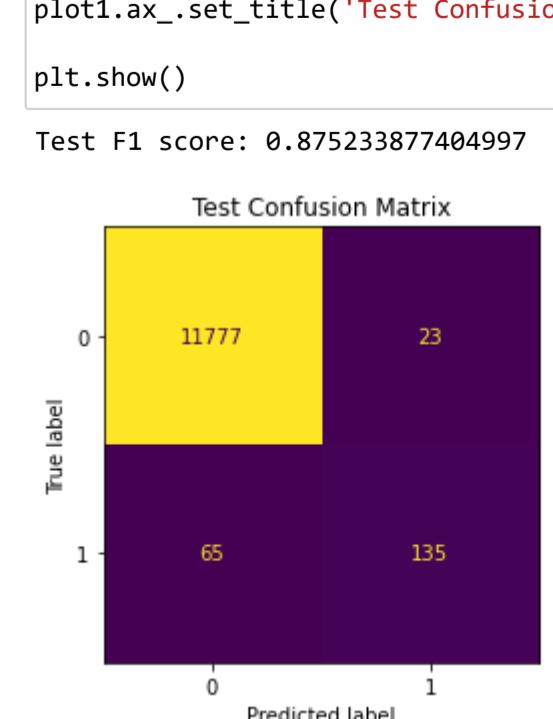
```
Random Search CV F1 score result
```



```
In [452]: x_cfl = LogisticRegression(C=0.1)
x_cfl.fit(df, y_d2.values.ravel())
Out[452]: LogisticRegression(C=0.1)
```

```
In [453]: test_df = test_dataset(num_of_models, model_list, X_ensemble_test)
y_test_pred = x_cfl.predict(test_df)
print("Test F1 score:", f1_score(y_ensemble_test, y_test_pred, average='macro'))

plot1 = plot_confusion_matrix(x_cfl, test_df, y_ensemble_test)
plot1.ax_.set_title('Test Confusion Matrix')
plot1.show()
Test F1 score: 0.875233877404997
```



Results

- 10 DT Model - 0.8706245599369046
- 20 DT Model - 0.8787366046249295
- 50 DT Model - 0.8817634631713474
- 75 DT Model - 0.875233877404997

Median Imputed Data

```
In [454]: X_train, X_ensemble_test, y_train, y_ensemble_test = train_test_split(data_df, y_true, stratify=y_true, test_size=0.2)
dataset_d1, dataset_d2, y_d1, y_d2 = train_test_split(X_train, y_train, stratify=y_train, test_size=0.5)
print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)

(24000, 149) (24000, 149) (12000, 149)

In [455]: high_correlated_columns = ['sensor45_measure', 'sensor32_measure', 'sensor46_measure', 'sensor47_measure', 'sensor67_measure']
for column in high_correlated_columns:
    if column in X_train.columns:
        dataset_d1 = dataset_d1.drop([column],axis=1)
        dataset_d2 = dataset_d2.drop([column],axis=1)
        X_ensemble_test = X_ensemble_test.drop([column],axis=1)

dataset_d1 = dataset_d1.drop(['sensor64_histogram_bin5'],axis=1)
dataset_d2 = dataset_d2.drop(['sensor64_histogram_bin5'],axis=1)
X_ensemble_test = X_ensemble_test.drop(['sensor64_histogram_bin5'],axis=1)

In [456]: median_values = dataset_d1.median()
columns = dataset_d1.columns
dataset_d1 = dataset_d1.fillna(median_values)
dataset_d2 = dataset_d2.fillna(median_values)
X_ensemble_test = X_ensemble_test.fillna(median_values)

print(dataset_d1.shape, dataset_d2.shape, X_ensemble_test.shape)

dataset_d1 = pd.DataFrame(dataset_d1)
dataset_d2 = pd.DataFrame(dataset_d2)
X_ensemble_test = pd.DataFrame(X_ensemble_test)

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns

(24000, 143) (24000, 143) (12000, 143)

In [457]: column_names = dict()
for column in dataset_d1.columns:
    value = column.split("_")
    if value[1] == "histogram":
        if value[0] not in column_names:
            column_names[value[0]] = []
        column_names[value[0]].append(column)

print(column_names)

{'sensor': ['sensor7_histogram_bin2', 'sensor7_histogram_bin3', 'sensor7_histogram_bin4', 'sensor7_histogram_bin5', 'sensor7_histogram_bin6', 'sensor7_histogram_bin7', 'sensor7_histogram_bin8', 'sensor7_histogram_bin9'], 'sensor24': ['sensor24_histogram_bin5', 'sensor24_histogram_bin6', 'sensor24_histogram_bin7', 'sensor24_histogram_bin8'], 'sensor25': ['sensor25_histogram_bin1', 'sensor25_histogram_bin2', 'sensor25_histogram_bin3', 'sensor25_histogram_bin4', 'sensor25_histogram_bin5', 'sensor25_histogram_bin6', 'sensor25_histogram_bin7', 'sensor25_histogram_bin8'], 'sensor26': ['sensor26_histogram_bin0', 'sensor26_histogram_bin1', 'sensor26_histogram_bin2', 'sensor26_histogram_bin3', 'sensor26_histogram_bin4', 'sensor26_histogram_bin5', 'sensor26_histogram_bin6', 'sensor26_histogram_bin7', 'sensor26_histogram_bin8'], 'sensor64': ['sensor64_histogram_bin1', 'sensor64_histogram_bin2', 'sensor64_histogram_bin3', 'sensor64_histogram_bin4', 'sensor64_histogram_bin5', 'sensor64_histogram_bin6', 'sensor64_histogram_bin7', 'sensor64_histogram_bin8'], 'sensor69': ['sensor69_histogram_bin1', 'sensor69_histogram_bin2', 'sensor69_histogram_bin3', 'sensor69_histogram_bin4', 'sensor69_histogram_bin5', 'sensor69_histogram_bin6', 'sensor69_histogram_bin7', 'sensor69_histogram_bin8'], 'sensor105': ['sensor105_histogram_bin0', 'sensor105_histogram_bin1', 'sensor105_histogram_bin2', 'sensor105_histogram_bin3', 'sensor105_histogram_bin4', 'sensor105_histogram_bin5', 'sensor105_histogram_bin6', 'sensor105_histogram_bin7', 'sensor105_histogram_bin8']}

In [458]: dataset_d1 = add_average_bin_feature(dataset_d1, column_names)
dataset_d2 = add_average_bin_feature(dataset_d2, column_names)
X_ensemble_test = add_average_bin_feature(X_ensemble_test, column_names)

In [459]: mnmx_scaler = MinMaxScaler()
columns = dataset_d1.columns

dataset_d1 = pd.DataFrame(mnmx_scaler.fit_transform(dataset_d1))
dataset_d2 = pd.DataFrame(mnmx_scaler.transform(dataset_d2))
X_ensemble_test = pd.DataFrame(mnmx_scaler.transform(X_ensemble_test))

dataset_d1.columns = columns
dataset_d2.columns = columns
X_ensemble_test.columns = columns

In [508]: # Number of models here is a hyper parameter and has been tried with multiple values and the results are below
num_of_models = 75
model_list, df = ensemble_model_DT(num_of_models, dataset_d1, y_d1, dataset_d2, y_d2)

Generating Samples Done
Fitting Base Models Done
Predict on Dataset 2 Done
Formed Meta Dataset Done

In [509]: # Logistic Regression as the meta classifier and it is hyper parameterized with the GridSearch CV
x_cfl = LogisticRegression()

params={
    'C':[10 ** x for x in range(-5, 4)]}
grid_cfl=GridSearchCV(x_cfl, param_grid=params, n_jobs=-1, scoring='f1')
grid_cfl.fit(df, y_d2.values.ravel())

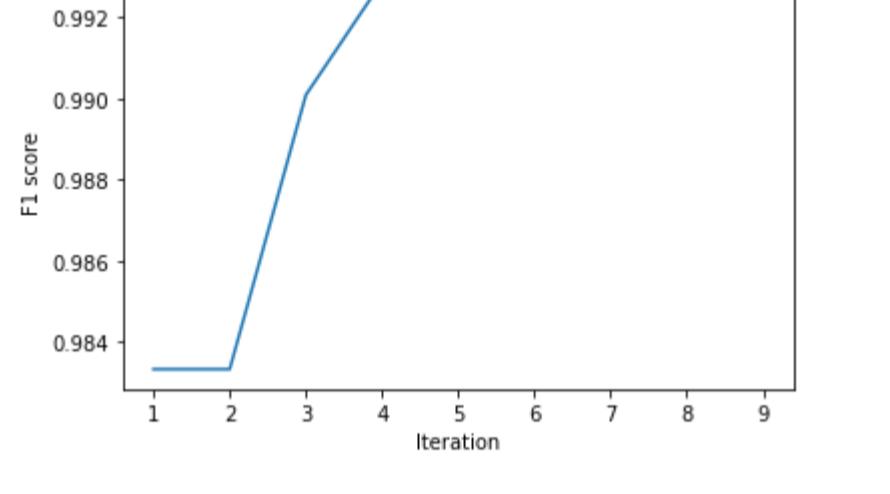
Out[509]: GridSearchCV(estimator=LogisticRegression(), n_jobs=-1,
param_grid={'C': [1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,
1000]})

In [510]: # plotting the various mean f1 scores
labels = grid_cfl.cv_results_['params']
x_axis = range(1,10)
y_axis = grid_cfl.cv_results_['mean_test_score']
for i, label in enumerate(labels):
    print(label,":", y_axis[i])

plt.xlabel('Iteration')
plt.ylabel('F1 score')
plt.title('Random Search CV F1 score result')
plt.plot(x_axis, y_axis)

print("BEST MODEL -- ", grid_cfl.best_estimator_)

{'C': 1e-05} : 0.9833333333333332
{'C': 0.0001} : 0.9833333333333332
{'C': 0.001} : 0.9900833333333333
{'C': 0.01} : 0.992875
{'C': 0.1} : 0.9930833333333334
{'C': 1} : 0.9926250000000001
{'C': 10} : 0.9925416666666667
{'C': 100} : 0.9925833333333334
{'C': 1000} : 0.9925833333333334
BEST MODEL -- LogisticRegression(C=0.1)

Random Search CV F1 score result


```
F1 score
iteration
```



| Iteration | F1 score |
|-----------|----------|
| 1         | 0.983    |
| 2         | 0.988    |
| 3         | 0.992    |
| 4         | 0.992    |
| 5         | 0.992    |
| 6         | 0.992    |
| 7         | 0.992    |
| 8         | 0.992    |
| 9         | 0.992    |

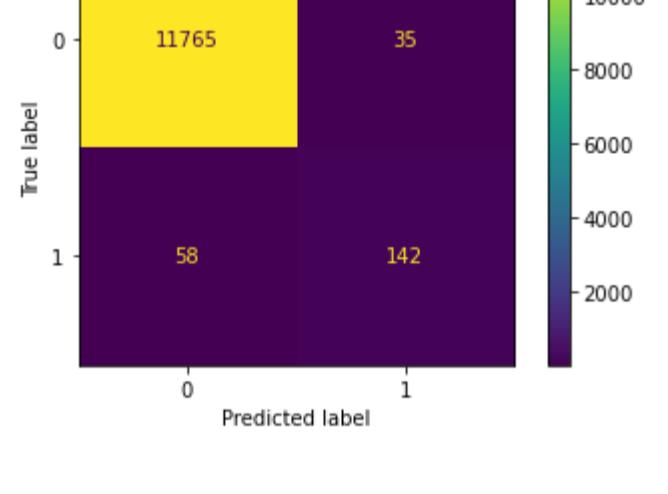


In [519]: x_cfl = LogisticRegression(C=0.1)
x_cfl.fit(df,y_d2.values.ravel())

Out[519]: LogisticRegression(C=0.1)

In [520]: test_df = test_dataset(num_of_models, model_list, X_ensemble_test)
y_test_pred = x_cfl.predict(test_df)
print("Test F1 score:", f1_score(y_ensemble_test, y_test_pred, average='macro'))

plot1 = plot_confusion_matrix(x_cfl, test_df, y_ensemble_test)
plot1.ax_.set_title('Test Confusion Matrix')
plt.show()

Test F1 score: 0.8746894043266515


| Predicted label |       | True label |
|-----------------|-------|------------|
| 0               | 1     | 0          |
| 0               | 11765 | 35         |
| 1               | 58    | 142        |



```
Test Confusion Matrix
```



| Predicted label | True label | Count |
|-----------------|------------|-------|
| 0               | 0          | 11765 |
| 0               | 1          | 35    |
| 1               | 0          | 58    |
| 1               | 1          | 142   |


```

The Ensemble Models with FF dataset

- ## Take aways from Modelling

```
In [548]: GREEN = '\033[92m'
END = '\033[0m'
BOLD = '\033[1m'

print(BOLD+"The Highlighted Rows are the ones with higher F1 score."+END)

table = PrettyTable()
table.field_names = ["Model", "Parameter", "Imputation", "Train F1 Score", "CV F1 score", "Test F1 Score"]
table.add_row(["KNN", "manhattan, n=3, distance", "Median", 1.0, 0.839, 0.847])
table.add_row(["KNN", "manhattan, n=3, distance", "Median", 1.0, 0.845, 0.831])
table.add_row(["KNN", "manhattan, n=4, distance", "KNN imputer", 1.0, 0.845, 0.831])
table.add_row(["XG Boost", "", ""])
table.add_row([BOLD+"Random Forest"+END, "d=100, n=1500", "Median", 1.0, 0.988, BOLD+GREEN+"0.889"+END])
table.add_row([BOLD+"XG Boost"+END, "ss=1, d=5, n=2000, lr=0.05, colsm_tree=0.5", "Median", 1.0, 0.913, BOLD+GREEN+"0.887"+END])
table.add_row([BOLD+"XG Boost"+END, "ss=0.5, d=10, n=2000, lr=0.2, colsm_tree=0.3", "KNN imputer", 1.0, 0.914, 0.877])
table.add_row(["Stacking Classifier", "RF, XGB, KNN, meta=KNN", "Median", 1.0, 0.915, BOLD+GREEN+"0.894"+END])
table.add_row(["Stacking Classifier", "RF, XGB, KNN, meta=KNN", "KNN imputer", 1.0, 0.906, 0.891])
table.add_row([BOLD+"Voting Classifier"+END, "RF, XGB, KNN", "Median", 1.0, 0.916, BOLD+GREEN+"0.898"+END])
table.add_row([BOLD+"Voting Classifier"+END, "RF, XGB, KNN", "KNN imputer", 1.0, 0.904, 0.889])
table.add_row(["Ensemble Models without FE", "10 DT + Logistic Regression", "Median", "0.99", "-", 0.839])
table.add_row(["Ensemble Models without FE", "20 DT + Logistic Regression", "Median", "0.99", "-", 0.833])
table.add_row(["Ensemble Models without FE", "50 DT + Logistic Regression", "Median", "0.99", "-", 0.877])
table.add_row(["Ensemble Models without FE", "75 DT + Logistic Regression", "Median", "0.99", "-", 0.876])
table.add_row(["Ensemble Models with FE", "10 DT + Logistic Regression", "Median", "0.99", "-", 0.843])
table.add_row(["Ensemble Models with FE", "20 DT + Logistic Regression", "Median", "0.99", "-", 0.874])
table.add_row([BOLD+"Ensemble Models with FE"+END, "50 DT + Logistic Regression", "Median", "0.99", "-", BOLD+GREEN+"0.881"+END])
table.add_row(["Ensemble Models with FE", "75 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.875])
table.add_row(["Ensemble Models with FE", "10 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.839])
table.add_row(["Ensemble Models with FE", "20 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.858])
table.add_row(["Ensemble Models with FE", "50 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.877])
table.add_row(["Ensemble Models with FE", "75 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.869])
table.add_row(["Ensemble Models without FE", "10 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.853])
table.add_row(["Ensemble Models without FE", "20 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.853])
table.add_row(["Ensemble Models without FE", "50 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.863])
table.add_row(["Ensemble Models without FE", "75 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.876])
table.add_row(["Ensemble Models with FE", "10 DT + Logistic Regression", "Median", "0.99", "-", 0.843])
table.add_row(["Ensemble Models with FE", "20 DT + Logistic Regression", "Median", "0.99", "-", 0.877])
table.add_row(["Ensemble Models with FE", "50 DT + Logistic Regression", "Median", "0.99", "-", 0.881])
table.add_row(["Ensemble Models with FE", "75 DT + Logistic Regression", "Median", "0.99", "-", 0.874])
table.add_row(["Ensemble Models with FE", "10 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.871])
table.add_row(["Ensemble Models with FE", "20 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.878])
table.add_row(["Ensemble Models with FE", "50 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.881])
table.add_row(["Ensemble Models with FE", "75 DT + Logistic Regression", "KNN imputer", "0.99", "-", 0.875])
```

```
table._max_width = {"Model":30, "Parameter":30, "Imputation":5, "CV F1 score":5, "Test F1 Score":5}
print(table)
```

The Highlighted Rows are the ones with higher F1 score.

Model	Parameter	Imputation	Train F1 Score	CV F1 score	Test F1 Score
KNN	manhattan, n=3, distance	Median	1.0	0.839	0.847
KNN	manhattan, n=4, distance	KNN imputer	1.0	0.845	0.831
Random Forest	d=100, n=1500	Median	1.0	0.988	BOLD+GREEN+"0.889"
Random Forest	d=100, n=1000	KNN imputer	1.0	0.983	0.879
XG Boost	ss=1, d=5, n=2000, lr=0.05, colsm_tree=0.5	Median	1.0	0.913	0.887
XG Boost	ss=0.5, d=10, n=2000, lr=0.2, colsm_tree=0.3	KNN imputer	1.0	0.914	0.877
Stacking Classifier	RF, XGB, KNN, meta=KNN	Median	1.0	0.915	BOLD+GREEN+"0.894"
Stacking Classifier	RF, XGB, KNN, meta=KNN	KNN imputer	1.0	0.906	0.891
Voting Classifier	RF, XGB, KNN	Median	1.0	0.916	BOLD+GREEN+"0.898"
Voting Classifier	RF, XGB, KNN	KNN imputer	1.0	0.904	0.889
Ensemble Models without FE	10 DT + Logistic Regression	Median	"0.99", "-," 0.839		
Ensemble Models without FE	20 DT + Logistic Regression	Median	"0.99", "-," 0.833		
Ensemble Models without FE	50 DT + Logistic Regression	Median	"0.99", "-," 0.877		
Ensemble Models without FE	75 DT + Logistic Regression	Median	"0.99", "-," 0.876		
Ensemble Models without FE	10 DT + Logistic Regression	KNN imputer	0.99	-	0.853
Ensemble Models without FE	20 DT + Logistic Regression	KNN imputer	0.99	-	0.853
Ensemble Models without FE	50 DT + Logistic Regression	KNN imputer	0.99	-	0.863
Ensemble Models without FE	75 DT + Logistic Regression	KNN imputer	0.99	-	0.876
Ensemble Models with FE	10 DT + Logistic Regression	Median	0.99	-	0.843
Ensemble Models with FE	20 DT + Logistic Regression	Median	0.99	-	0.877
Ensemble Models with FE	50 DT + Logistic Regression	Median	0.99	-	0.881
Ensemble Models with FE	75 DT + Logistic Regression	Median	0.99	-	0.874
Ensemble Models with FE	10 DT + Logistic Regression	KNN imputer	0.99	-	0.871
Ensemble Models with FE	20 DT + Logistic Regression	KNN imputer	0.99	-	0.878
Ensemble Models with FE	50 DT + Logistic Regression	KNN imputer	0.99	-	0.881
Ensemble Models with FE	75 DT + Logistic Regression	KNN imputer	0.99	-	0.875

Deployment

The solution has been deployed on the following website

- <http://ec2-52-15-225-42.us-east-2.compute.amazonaws.com:8080/index> (<http://ec2-52-15-225-42.us-east-2.compute.amazonaws.com:8080/index>)