

# Notebook Index

## 1- Import essential libraries

## 2- Read the data and do exploratory data analysis

### 2.1 Check for duplicates

### 2.2 Check data types and null values

### 2.3 Check for classes balance

### 2.3 result: There is imbalance, we will deal with if the score are affected

### 2.4 Merge the data with time as Date/Time and drop the columns Data and Time

### 2.5 Check features correlation

### 2.6 Aggregating features to deal with the correlation issue

### 2.7 Check the correlation after aggregation

### 2.8 Pair plot after the aggregation

### 3- Define and prepare the features and target

### 3.1 Split the dataset into train and test data

### 3.2 Scale the features (Standardization)

### 4- Initiate the machine learning models

### 4.1 Logistic Regression Model with evaluation metrics and coefficients

### 4.1 Result: The bench mark model. it has a great overall scores. but the individual classes (2 and 3) has a low scores

### 4.2 Gaussian naïve bayes with evaluation metrics and class probabilities

### 4.2 Result: The second best classifier out of the three classifiers. The score of individual classes are good except Class (3) has a low scores but still acceptable

### 4.3 XGboost with evaluation metrics and features importance plot

### 4.3 Result: The highest score out of the three models. Also the score for the individual classes are high

## 1- Import essential libraries

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# 2- Read the data and do exploratory data analysis

df = pd.read_csv('Occupancy_Estimation.csv')
df.head()
```

	Date	Time	S1_Temp	S2_Temp	S3_Temp	S4_Temp	S1_Light	S2_Light	S3_Light	S4_Light	S1_Sound	S2_Sound	S3_Sound	S4_Sound	SS_CO2	SS_CO2_Slope	S6_PIR	S7_PIR	Room	Occupancy_Count
0	2017/12/22	10:40:41	24.94	24.75	24.56	25.44	121	34	53	40	0.98	0.19	0.06	0.06	390	0.769231	0	0		0
1	2017/12/22	10:50:12	24.94	24.74	24.56	25.44	121	33	53	40	0.98	0.05	0.06	0.06	390	0.646154	0	0		0
2	2017/12/22	10:50:42	25.00	24.75	24.50	25.44	121	34	53	40	0.43	0.11	0.06	0.06	390	0.519231	0	0		0
3	2017/12/22	10:51:13	25.00	24.74	24.56	25.44	121	34	53	40	0.43	0.10	0.10	0.09	390	0.388462	0	0		0
4	2017/12/22	10:51:44	25.00	24.75	24.56	25.44	121	34	54	40	0.18	0.06	0.06	0.06	390	0.253846	0	0		0

```
In [3]: df.shape
Out[3]: (10129, 19)

# 2.1 Check for duplicates
df[df.duplicated()]

# 2.2 Check data types and null values
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10129 entries, 0 to 10128
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  --
0   Date                   10129 non-null object
1   Time                   10129 non-null object
2   S1_Temp                10129 non-null float64
3   S2_Temp                10129 non-null float64
4   S3_Temp                10129 non-null float64
5   S4_Temp                10129 non-null float64
6   S1_Light               10129 non-null int64
7   S2_Light               10129 non-null int64
8   S3_Light               10129 non-null int64
9   S4_Light               10129 non-null int64
10  S1_Sound               10129 non-null float64
11  S2_Sound               10129 non-null float64
12  S3_Sound               10129 non-null float64
13  S4_Sound               10129 non-null float64
14  SS_CO2                 10129 non-null float64
15  SS_CO2_Slope           10129 non-null float64
16  S6_PIR                 10129 non-null int64
17  S7_PIR                 10129 non-null int64
18  Room_Occupancy_Count   10129 non-null int64
dtypes: float64(9), int64(8), object(2)
memory usage: 1.5+ MB
```

```
In [7]: df.describe()
Out[7]:
```

	S1_Temp	S2_Temp	S3_Temp	S4_Temp	S1_Light	S2_Light	S3_Light	S4_Light	S1_Sound	S2_Sound	S3_Sound	S4_Sound	SS_CO2	SS_CO2_Slope	S6_PIR	S7_PIR	Room	Occupancy_Count
count	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000
mean	25.454012	25.500659	25.056211	25.754125	25.445059	26.05629	34.248484	13.220259	0.168178	0.120066	0.116119	0.103840	0.406040	0.137391	0.430633	0.474077		0
std	0.350151	0.506325	0.427283	0.356434	0.111254	0.730417	56.400744	19.602219	0.316709	0.269503	0.413637	0.120683	0.19946490	0.706600	0.550500	0.3450000		0
min	24.940000	24.700000	24.400000	24.400000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000		0
25%	25.100000	25.100000	24.900000	25.400000	0.000000	0.000000	0.000000	0.000000	0.070000	0.050000	0.040000	0.040000	0.000000	0.000000	0.000000	0.000000		0
50%	25.300000	25.300000	24.900000	25.700000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000		0
75%	25.600000	25.600000	25.300000	26.000000	0.000000	0.000000	14.000000	50.000000	0.220000	0.080000	0.060000	0.070000	0.000000	0.000000	0.000000	0.000000		0
max	26.300000	29.000000	26.100000	26.500000	105.000000	295.000000	280.000000	74.000000	3.880000	3.440000	3.670000	3.400000	3.400000	3.400000	3.400000	3.400000		0

### 2.3 Check for classes balance

```
In [8]: df.Room_Occupancy_Count.value_counts()
Out[8]:
0      8228
1       748
2       390
3       459
Name: Room_Occupancy_Count, dtype: int64
```

### 2.3 result: There is imbalance, we will deal with if the score are affected

### 2.4 Merge the data with time as Date/Time and drop the columns Data and Time

```
In [9]: dateTimeConcat = df[['date','time']] + ' ' + df[['time']]
df[dateTimeConcat] = pd.to_datetime(dateTimeConcat)

In [10]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10129 entries, 0 to 10128
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  --
0   Date                   10129 non-null object
1   Time                   10129 non-null object
2   S1_Temp                10129 non-null float64
3   S2_Temp                10129 non-null float64
4   S3_Temp                10129 non-null float64
5   S4_Temp                10129 non-null float64
6   S1_Light               10129 non-null int64
7   S2_Light               10129 non-null int64
8   S3_Light               10129 non-null int64
9   S4_Light               10129 non-null int64
10  S1_Sound               10129 non-null float64
11  S2_Sound               10129 non-null float64
12  S3_Sound               10129 non-null float64
13  S4_Sound               10129 non-null float64
14  SS_CO2                 10129 non-null float64
15  SS_CO2_Slope           10129 non-null float64
16  S6_PIR                 10129 non-null int64
17  S7_PIR                 10129 non-null int64
18  Room_Occupancy_Count   10129 non-null int64
19  dateTimeConcat         10129 non-null datetime64[ns]
dtypes: datetime64[ns](1), float64(9), int64(8), object(2)
memory usage: 1.5+ MB
```

```
In [11]: df.drop(columns=['date','time'], inplace = True)

In [12]: df.columns
Out[12]:
Index(['S1_Temp', 'S2_Temp', 'S3_Temp', 'S4_Temp', 'S1_Light', 'S2_Light',
       'S3_Light', 'S4_Light', 'S1_Sound', 'S2_Sound', 'S3_Sound', 'S4_Sound',
       'SS_CO2', 'SS_CO2_Slope', 'S6_PIR', 'S7_PIR', 'Room_Occupancy_Count',
       'dateTimeConcat'],
      dtype='object')
```

### 2.5 Check features correlation

```
In [13]: df.corr()
Out[13]:
```

	S1_Temp	S2_Temp	S3_Temp	S4_Temp	S1_Light	S2_Light	S3_Light	S4_Light	S1_Sound	S2_Sound	S3_Sound	S4_Sound	SS_CO2	SS_CO2_Slope	S6_PIR	S7_PIR
S1_Temp	1.000000	0.799707	0.948839	0.858270	0.686743	0.548735	0.645163	0.212217	0.436009	0.391137	0.438769	0.355826	0.666718	0.137391	0.430633	0.474077
S2_Temp	0.799707	1.000000	0.765255	0.695681	0.639773	0.645987	0.607349	0.370897	0.482774	0.409545	0.427133	0.378724	0.743722	0.202547	0.470901	0.465884
S3_Temp	0.948839	0.765255	1.000000	0.883136	0.594311	0.500054	0.642001	0.301419	0.375183	0.344025	0.398117	0.320182	0.621308	0.095842	0.403395	0.460309
S4_Temp	0.858270	0.695681	0.883136	1.000000	0.595482	0.456350	0.588459	0.388971	0.395111	0.312394	0.340408	0.294839	0.693020	0.105208	0.340000	0.332007
S1_Light	0.686743	0.639773	0.594311	0.595482	1.000000	0.940090	0.810180	0.519385	0.801188	0.534774	0.484880	0.447172	0.602760	0.486185	0.407159	0.545213
S2_Light	0.548735	0.645987	0.500054	0.456350	0.940090	1.000000	0.705978	0.458904	0.503051	0.556030	0.430959	0.414902	0.566764	0.402051	0.554658	0.562797
S3_Light	0.645163	0.607349	0.642001	0.588459	0.810180	0.705978	1.000000	0.579484	0.502066	0.434659	0.577151	0.474066	0.596395	0.447708	0.501636	0.577615
S4_Light	0.212217	0.370897	0.301419	0.388971	0.519385	0.458904	0.579484	1.000000	0.293032	0.302948	0.169702	0.239793	0.346608	0.232178	0.324545	0.220196
S1_Sound	0.436009	0.482774	0.375183	0.395111	0.801188	0.503051	0.502066	0.293032	1.000000	0.560062	0.540736	0.557713	0.391303	0.318515	0.342015	0.463040
S2_Sound	0.391137	0.409545	0.344025	0.312394	0.534774	0.556030	0.430959	0.303949	0.560062	1.000000	0.529830	0.578035	0.333386	0.357235	0.485697	0.577231
S3_Sound	0.438769	0.427133	0.398117	0.340408	0.494800	0.430959	0.577151	0.169702	0.540736	0.529830	1.000000	0.696670	0.447220	0.318515	0.342535	0.528620
S4_Sound	0.355826	0.743722	0.320182	0.294839	0.447172	0.413932	0.473606	0.200793	0.557733	0.578635	0.696670	1.000000	0.306620	0.328519	0.394954	0.466848
SS_CO2	0.666718	0.743722	0.621308	0.693020	0.602740	0.566764	0.650829	0.149608	0.391903	0.333636	0.447220	0.330629	1.000000	0.069220	0.399626	0.473437
SS_CO2_Slope	0.137391	0.202547	0.095842	0.105208	0.486185	0.402051	0.447708	0.212718	0.335772	0.357235	0.318515	0.325519	0.069220	1.000000	0.366374	0.423546
S6_PIR	0.474077	0.465884	0.460309	0.332007	0.607159	0.546568	0.501636	0.324545	0.622015	0.485697	0.434225	0.394954	0.577125	0.366374	1.000000	0.571125
S7_PIR	0.474077	0.465884	0.460309	0.332007	0.546213	0.556797	0.577615	0.222195	0.483040	0.507231	0.536800	0.466848	0.741367	0.425346	0.571125	1.000000
Room_Occupancy_Count	0.700886	0.671263	0.652047	0.526595	0.480558	0.788754	0.793081	0.355715	0.573748	0.557783	0.536105	0.460287	0.601144	0.601105	0.633133	0.695138

### 2.6 Aggregating features to deal with the correlation issue

```
In [14]: df[['S1_4_Temp_Mean']] = (df[['S1_Temp']] + df[['S2_Temp']] + df[['S3_Temp']] + df[['S4_Temp']]) / 4
df[['S1_4_Light_Mean']] = (df[['S1_Light']] + df[['S2_Light']] + df[['S3_Light']] + df[['S4_Light']]) / 4
df[['S1_4_Sound_Mean']] = (df[['S1_Sound']] + df[['S2_Sound']] + df[['S3_Sound']] + df[['S4_Sound']]) / 4
df[['S6_7_PIR_Sum']] = df[['S6_PIR']] + df[['S7_PIR']]

In [15]: #Original_Features = ['S1_Temp', 'S2_Temp', 'S3_Temp', 'S4_Temp', 'S1_Light', 'S2_Light',
#                        'S3_Light', 'S4_Light', 'S1_Sound', 'S2_Sound', 'S3_Sound', 'S4_Sound',
#                        'SS_CO2', 'SS_CO2_Slope', 'S6_PIR', 'S7_PIR']
#Aggregated_Features = ['S1_4_Temp_Mean', 'S1_4_Light_Mean', 'S1_4_Sound_Mean',
#                        'S6_7_PIR_Sum', 'SS_CO2', 'SS_CO2_Slope']

In [16]: list_df_Agg = Aggregated_Features.copy()
list_df_Agg.insert(0, 'Room_Occupancy_Count')

In [17]: df[list_df_Agg].describe()
```

	Room_Occupancy_Count	S1_4_Temp_Mean	S1_4_Light_Mean	S1_4_Sound_Mean	S6_7_PIR_Sum	SS_CO2	SS_CO2_Slope
count	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000	10129.000000
mean	0.398559	25.452704	24.732525	0.137551	0.597711	0.406040	0.137391
std	0.896333	0.399201	43.11460	0.231763	0.493767	0.19946490	0.154990
min	0.000000	24.800000	0.000000	0.000000	0.000000	0.3450000	0.290154
25%	0.000000	25.127000	0.000000	0.025500	0.000000	0.3550000	0.040154
50%	0.000000	25.300000	0.000000	0.067000	0.000000	0.3600000	0.000000
75%	0.000000	25.627000	0.000000	0.079000	0.000000	0.4500000	0.000000
max	3.000000	26.800000	193.750000	3.182500	2.000000	1270.000000	0.890769

### 2.7 Check the correlation after aggregation

```
In [18]: df[list_df_Agg].corr()
Out[18]:
```

	Room_Occupancy_Count	S1_4_Temp_Mean	S1_4_Light_Mean	S1_4_Sound_Mean	S6_7_PIR_Sum	SS_CO2	SS_CO2_Slope
Room_Occupancy_Count	1.000000	0.692698	0.854160	0.653527	0.748047	0.660144	0.601105
S1_4_Temp_Mean	0.692698	1.000000	0.685673	0.500396	0.523816	0.637725	0.151397
S1_4_Light_Mean	0.854160	0.685673	1.000000	0.646892	0.667364	0.626629	0.507454
S1_4_Sound_Mean	0.653527	0.500396	0.646892	1.000000	0.657848	0.472438	0.401771
S6_7_PIR_Sum	0.748047	0.523816	0.667364	0.657848	1.000000	0.488761	0.448805
SS_CO2	0.660144	0.637725	0.626629	0.472438	0.488761	1.000000	0.069220
SS_CO2_Slope	0.601105	0.151397	0.507454	0.401771	0.448805	0.069220	1.000000

### 2.8 Pair plot after the aggregation

```
In [19]: sns.pairplot(df[list_df_Agg])
```