

```
In [1]: from IPython.display import Image, Latex  
Image(filename='title.jpg',width=1000)
```

Out[1]:



Modeling of CO Emissions from Cars

Christine Al-Thifairy (97074755)

Elisabeth Grasia Putri (611410005)

Kin Man Lam (15823898)

Assignment 2 for COMP7023 Predictive Analytics

School of Computer, Data and Mathematical Sciences,

Western Sydney University

Spring, 2022

```
In [2]: Image(filename='declar.jpg',width=900)
```

Out[2]:

Declaration

By including this statement, we the authors of this work, verify that:

- We hold a copy of this assignment that we can produce if the original is lost or damaged.
- We hereby certify that no part of this assignment/product has been copied from any other student's work or from any other source except where due acknowledgement is made in the assignment.
- No part of this assignment/product has been written/produced for us by another person except where such collaboration has been authorised by the subject lecturer/tutor concerned.
- We are aware that this work may be reproduced and submitted to plagiarism detection software programs for the purpose of detecting possible plagiarism (which may retain a copy on its database for future plagiarism checking).
- We hereby certify that we have read and understand what the School of Computer, Data and Mathematical Sciences defines as minor and substantial breaches of misconduct as outlined in the learning guide for this unit.

Modeling of CO Emissions From Cars

1 The Dataset

This assignment was developed based on the dataset provided by [Vehicle Certification Agency / UK](#) which includes 41 car manufacturers.

The dataset is uploaded to our group GitHub repo under the file name "Euro_6_latest_07-10-2022.zip".

Most of data are based on the new European driving cycle, Worldwide harmonized Light vehicles Test Procedure (WLTP).

The data is applied to the analysis of emissions including CO, THC, NOx, THC + NO, Particulates, and Noise emissions. Our focus in this report will be on CO emissions.

You can access this file on GitHub by clicking on this link: [assignment-2-pa-22-group-5.ipynb](#)

2 Evaluation Metrics

We will evaluate the performance of the models using the following metrics:

- MSE (Mean Square Error) to evaluate regression models.
- Confusion Matrix, and Accuracy for classification models.

3 Data Preprocessing

3.1 - Data Cleaning

First we import Python libraries and read the dataset.

```
In [3]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.model_selection import KFold

# read dataset file

df = pd.read_csv('Euro_6_latest.csv', delimiter=',', encoding='ISO-8859-1')
```

First, let's check the dim of the dataset.

```
In [4]: print('Dataset dim is ',df.shape)
```

Dataset dim is (4657, 44)

```
In [5]: %%html
<style>
  xtable {margin-left: 0 !important;}
  table.a td, table.a th {text-align: left !important}
  table.a th:first-child, table.a td:first-child {width: 25%}
</style>
```

Initial cleaning includes remove unwanted columns including:

1- Drop any column with more than 70% missing values

2- Drop the following unnecessary columns:

Columns	Reason
Transmission	This is a low-level category of car transmissions with 39 categories. We are going to use 'Manual or Automatic' which is a high-level category of vehicle transmissions with only 3 values.
Euro Standard Testing Scheme Date of change	Those columns are not relevant to prediction and won't help the model learning new trends
WLTP Imperial Low WLTP Imperial Medium WLTP Imperial High WLTP Imperial Extra High WLTP Imperial Combined WLTP Imperial Combined (Weighted)	Imperial measurements, we are using metric measurements.
Diesel VED Supplement	This is only relevant to Diesel cars
Emissions NOx [mg/km] THC Emissions [mg/km] Noise Level dB(A)	These features are targets (emissions). We are only predicting CO emission, therefore no need to keep the others.

```
In [6]: # Remove a funny column called Unnamed
df = df.loc[:, ~df.columns.str.contains('^Unnamed')]

# Remove features with more than 70% missing values
limitPer = len(df) * .70
df = df.dropna(thresh=limitPer, axis=1).copy()

# drop unnecessary columns
df = df.drop(['Transmission', 'Euro Standard', 'Testing Scheme', 'Date of change',
              'WLTP Imperial Medium', 'WLTP Imperial High', 'WLTP Imperial Extra',
              'WLTP Imperial Combined', 'WLTP Imperial Combined (Weighted)',
              'Diesel VED Supplement', 'Emissions NOx [mg/km]',
              'THC Emissions [mg/km]', 'Noise Level dB(A)'], axis=1).copy()

df.shape
```

Out[6]: (4657, 18)

Next, we give 'Manual or Automatic' column a new name 'Transmission' which is less verbos.

First, We look at the labels in 'Manual or Automatic' feature.

There are three labels in this column:

- 'Manual'
- 'Automatic'
- 'Electric - Not Applicable'

We going to create a new column called *Transmission* with all values from 'Manual or Automatic' column. However, we replace the label 'Electric - Not Applicable' with one word 'Electric' which is less verbos.

```
In [7]: # Create a new column 'Transmission', assign it a value of Automatic where the
# first character of Transmission is A
# or Manual if the first letter is M or Electric if first letter is E
AUTOMATIC = "Automatic"
MANUAL = "Manual"
ELECTRIC = "Electric"
df.loc[df['Manual or Automatic'].str.startswith('A'),'Transmission'] = AUTOMATIC
df.loc[df['Manual or Automatic'].str.startswith('M'),'Transmission'] = MANUAL
df.loc[df['Manual or Automatic'].str.startswith('E'),'Transmission'] = ELECTRIC
```

3.2 - Merge Low-Level Categories To Upper-Level Categories

To simplify the encoding process and reduce number of columns, we are going to merge few labels together in both *Fuel Type* and *Powertrain* columns.

First, let's examin the original labels in both columns:

```
In [8]: print('Fuel Type labels are:\n\n',df['Fuel Type'].unique(),'\n')
print('Powertrain labels are \n\n',df['Powertrain'].unique())
```

Fuel Type labels are:

```
['Petrol' 'Diesel' 'Electricity / Petrol' 'Petrol Electric' 'Electricity'
'Petrol / LPG' 'Diesel Electric' 'Electricity / Diesel']
```

Powertrain labels are

```
['Internal Combustion Engine (ICE)'
'Plug-in Hybrid Electric Vehicle (PHEV)'
'Mild Hybrid Electric Vehicle (MHEV)'
'Battery Electric Vehicle (BEV) / Pure Electric Vehicle / Electric Vehicle (EV)'
'Hybrid Electric Vehicle (HEV)' 'Micro Hybrid']
```

Fuel Type Column:

A new column called just *Fuel* to replace the original *Fuel Type* is added with only four labels **Pertrol, Diesel, Hybrid, and Electric**.

- Petrol, Petrol / LPG ---> Petrol
- Diesel ---> Diesel
- Electricity ---> Electric
- Electricity / Petrol, Petrol Electric, Diesel Electric, Electricity / Diesel ---> hybrid

Powertrain Column:

A new column called *PT* to replace the original Powertrain column with only 3 labels **ICE**, **EV**, and **Hybrid**

- Internal Combustion Engine (ICE) ---> ICE
- Plug-in Hybrid Electric Vehicle (PHEV), Mild Hybrid Electric Vehicle (MHEV), Hybrid Electric Vehicle (HEV), 'Micro Hybrid' ---> Hybrid
- Battery Electric Vehicle (BEV) / Pure Electric Vehicle / Electric Vehicle (EV) ---> EV

The new labels are less verbos and easy to encode.

```
In [9]: petrol = ['Petrol', 'Petrol / LPG']
diesel = ['Diesel']
electric = ['Electricity']
hybrid = ['Electricity / Petrol', 'Petrol Electric', 'Diesel Electric', 'Electric']
#df = df.copy()
df.loc[df['Fuel Type'].isin(petrol), 'Fuel'] = 'Petrol'
df.loc[df['Fuel Type'].isin(diesel), 'Fuel'] = 'Diesel'
df.loc[df['Fuel Type'].isin(electric), 'Fuel'] = 'Electric'
df.loc[df['Fuel Type'].isin(hybrid), 'Fuel'] = 'Hybrid'
```

```
In [10]: ice = ['Internal Combustion Engine (ICE)']
hybrid = ['Plug-in Hybrid Electric Vehicle (PHEV)',
          'Mild Hybrid Electric Vehicle (MHEV)',
          'Hybrid Electric Vehicle (HEV)', 'Micro Hybrid']
ev = ['Battery Electric Vehicle (BEV) / Pure Electric Vehicle / Electric Vehicle']

df.loc[df['Powertrain'].isin(ice), 'PT'] = 'ICE'
df.loc[df['Powertrain'].isin(ev), 'PT'] = 'EV'
df.loc[df['Powertrain'].isin(hybrid), 'PT'] = 'Hybrid'
```

And we drop the original columns Powertrain, Fuel Type, and Manual or Automatic.

```
In [11]: # Drop the old columns
df = df.drop(['Powertrain', 'Fuel Type', 'Manual or Automatic'], axis=1)
```

And rename PT back to Powertrain

```
In [12]: df.rename(columns={'PT': 'Powertrain'}, inplace=True)
```

3.3 - Handle Missing Values And Zeros

We check missing values in each of the remaining columns

```
In [13]: for (columnName, columnData) in df.iteritems():
          print(columnName, ' Num of empty cells : ', columnData.isnull().sum())
```

```
Manufacturer Num of empty cells : 0
Model Num of empty cells : 0
Description Num of empty cells : 0
Engine Capacity Num of empty cells : 2
Engine Power (PS) Num of empty cells : 221
Engine Power (Kw) Num of empty cells : 89
WLTP Metric Low Num of empty cells : 5
WLTP Metric Medium Num of empty cells : 5
WLTP Metric High Num of empty cells : 5
WLTP Metric Extra High Num of empty cells : 6
WLTP Metric Combined Num of empty cells : 15
WLTP Metric Combined (Weighted) Num of empty cells : 1000
WLTP CO2 Num of empty cells : 2
WLTP CO2 Weighted Num of empty cells : 1209
Emissions CO [mg/km] Num of empty cells : 108
Transmission Num of empty cells : 0
Fuel Num of empty cells : 0
Powertrain Num of empty cells : 0
```

Engine Capacity has two missing values which correspond to rows of Electric cars. Most likely because it is not applicable for Electric cars. So, we replace those missing values with zero.

```
In [14]: df['Engine Capacity'] = df['Engine Capacity'].fillna(0)
```

Engine Power (PS) column has 221 missing values. We will use ffill (forward filling) and bfill (backward filling) per group to fill the missing values of Engine Power (PS). We group by car manufacturer and Model. So, the function will search for another car from same manufacturer and same model and has value in that column and copies that value to the other one with missing value.

Engine Power (Kw) has 89 missing values. We apply the same method above to fill the missing values.

```
In [15]: df['Engine Power (PS)'] = df.groupby(['Manufacturer', 'Model'], sort=False) \
          ['Engine Power (PS)'].apply(lambda x: x.ffill().bfill())

# Same with Engine Power (Kw)
df['Engine Power (Kw)'] = df.groupby(['Manufacturer', 'Model'], sort=False) \
          ['Engine Power (Kw)'].apply(lambda x: x.ffill().bfill())
```

Some of the features has many zeros. For example, *WLTP Metric Combined (Weighted)* and *WLTP CO2*.

Let's count how many zeros in *WLTP Metric Combined (Weighted)* and *WLTP CO2 Weighted* columns

```
In [16]: print('WLTP Metric Combined (Weighted) zeros =', len(df) - np.count_nonzero(df['WLTP Metric Combined (Weighted)']))
          print('WLTP CO2 Weighted zeros =', len(df) - np.count_nonzero(df['WLTP CO2 Weighted']))
```

```
WLTP Metric Combined (Weighted) zeros = 3472
WLTP CO2 Weighted zeros = 3098
```

There are too many zeros values in both columns. So we are going to drop those 2 columns

```
In [17]: df = df.drop(['WLTP Metric Combined (Weighted)', 'WLTP CO2 Weighted'], axis=1)
```

For the rest of the features, we are going to fill missing values with mean of each column.

We did not replace 0 values with the mean because that is not correct, for example electric cars may have 0 CO emissions.

```
In [18]: df['WLTP Metric Combined'].fillna((df['WLTP Metric Combined'].mean()), inplace=True)
df['WLTP Metric Low'].fillna((df['WLTP Metric Low'].mean()), inplace=True)
df['WLTP Metric Medium'].fillna((df['WLTP Metric Medium'].mean()), inplace=True)
df['WLTP Metric High'].fillna((df['WLTP Metric High'].mean()), inplace=True)
df['WLTP Metric Extra High'].fillna((df['WLTP Metric Extra High'].mean()), inplace=True)
df['WLTP CO2'].fillna((df['WLTP CO2'].mean()), inplace=True)
df['Emissions CO [mg/km]'].fillna((df['Emissions CO [mg/km]'].mean()), inplace=True)
df.isnull().sum(axis = 0)
```

```
Out[18]: Manufacturer      0
Model                      0
Description                0
Engine Capacity            0
Engine Power (PS)         177
Engine Power (Kw)          89
WLTP Metric Low           0
WLTP Metric Medium        0
WLTP Metric High          0
WLTP Metric Extra High    0
WLTP Metric Combined      0
WLTP CO2                   0
Emissions CO [mg/km]      0
Transmission              0
Fuel                      0
Powertrain                 0
dtype: int64
```

There are still 177 missing values in Engine Power (PS) and 89 in Engine Power (Kw). We looked at the original dataset csv file. It turns out those car models have no data. Therefore we are going to drop those models.

```
In [19]: df = df.dropna().copy()
```

Finally, let's check the size of the dataset after all the preprocessing steps above.

```
In [20]: print('The dim of the dataset is ', df.shape, ' and original dim is (4657, 19)')
```

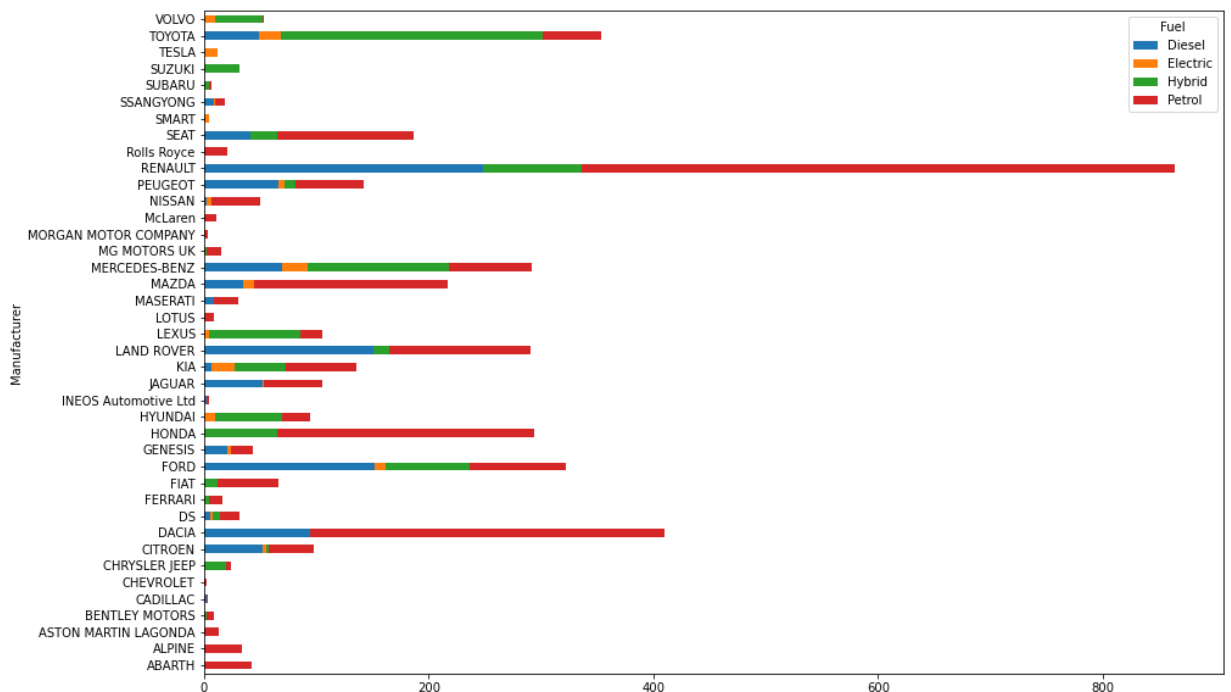
```
The dim of the dataset is (4467, 16) and original dim is (4657, 19)
```

4 Data Exploration

4.1 - Total Cars By Car Manufacturers

First we plot number of cars in the dataset by manufacturers and type of fuel. There are more models of some cars than others, especially the european models. It is obvious that Petrol still the dominant type of petrol. However, many car manufacturers are producing hybrid models.

```
In [21]: from IPython.core.display import display, HTML
display(HTML("<style>div.output_scroll { height: 44em; }</style>"))
plt.rcParams["figure.figsize"] = (15,10)
cols = ['Manufacturer', 'Fuel']
df1 = df[cols].copy()
df1.groupby(['Manufacturer', 'Fuel']).size().unstack().plot(sort_columns='Manu
kind='barh', stacked=True,width=0.5,linewidth=0.5);
```



4.2 - CO Emission By Car Manufacturer

Next, we plot the CO emission value against each car. Some of the cars show higher than others. Sometimes, this is because of number of models used for each cars is higher than other.

Fiat and Lotus are generating more CO emission than other cars in the dataset.

Note: Nissan CO Emission figures are recorded as zero in the original dataset.

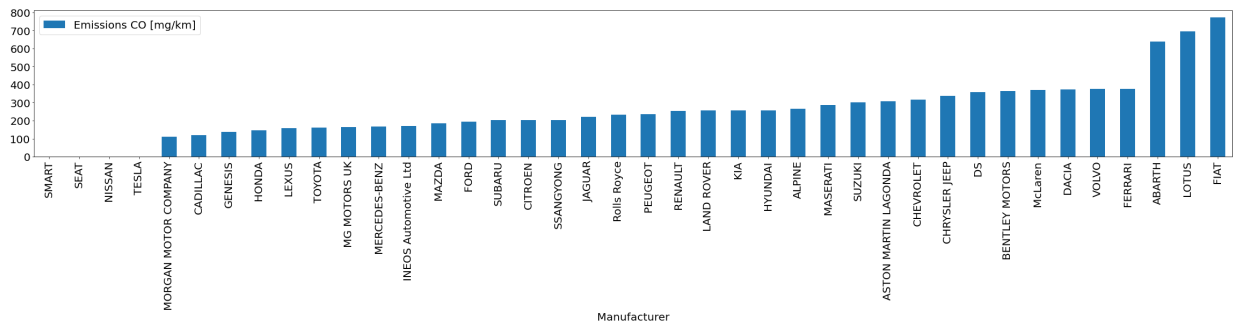
```
In [22]: cols = ['Manufacturer', 'Emissions CO [mg/km]']
plt.rcParams["figure.figsize"] = (40,5)

df2 = df[cols].copy()
df2 = df2.groupby('Manufacturer', as_index=False)['Emissions CO [mg/km]'].mean()
```



```
df2.sort_values('Emissions CO [mg/km]',inplace=True)
df2.plot(kind='bar',x='Manufacturer',y='Emissions CO [mg/km]',fontsize=20)
plt.legend(fontsize = 20)
plt.xlabel('Manufacturer', fontsize=20)
```

Out[22]: Text(0.5, 0, 'Manufacturer')

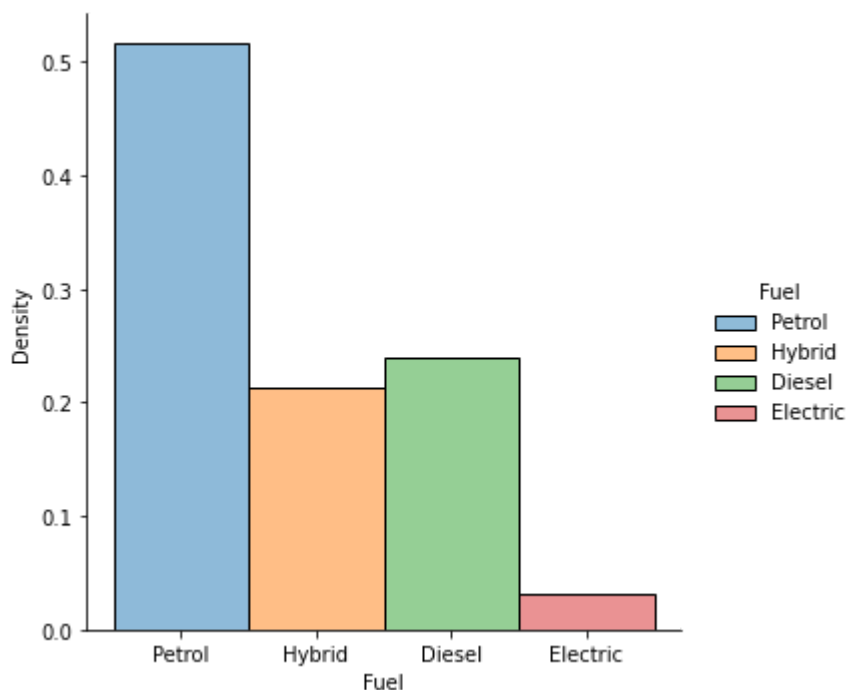


4.3 - Percentage Of Cars By Fuel Type

About half of the cars in the dataset are running on petrol.

```
In [23]: sns.displot(data=df1, x="Fuel",stat='density',kind='hist', hue=df['Fuel'])
```

Out[23]: <seaborn.axisgrid.FacetGrid at 0x7fec5dc530a0>



4.4 - Pairplot

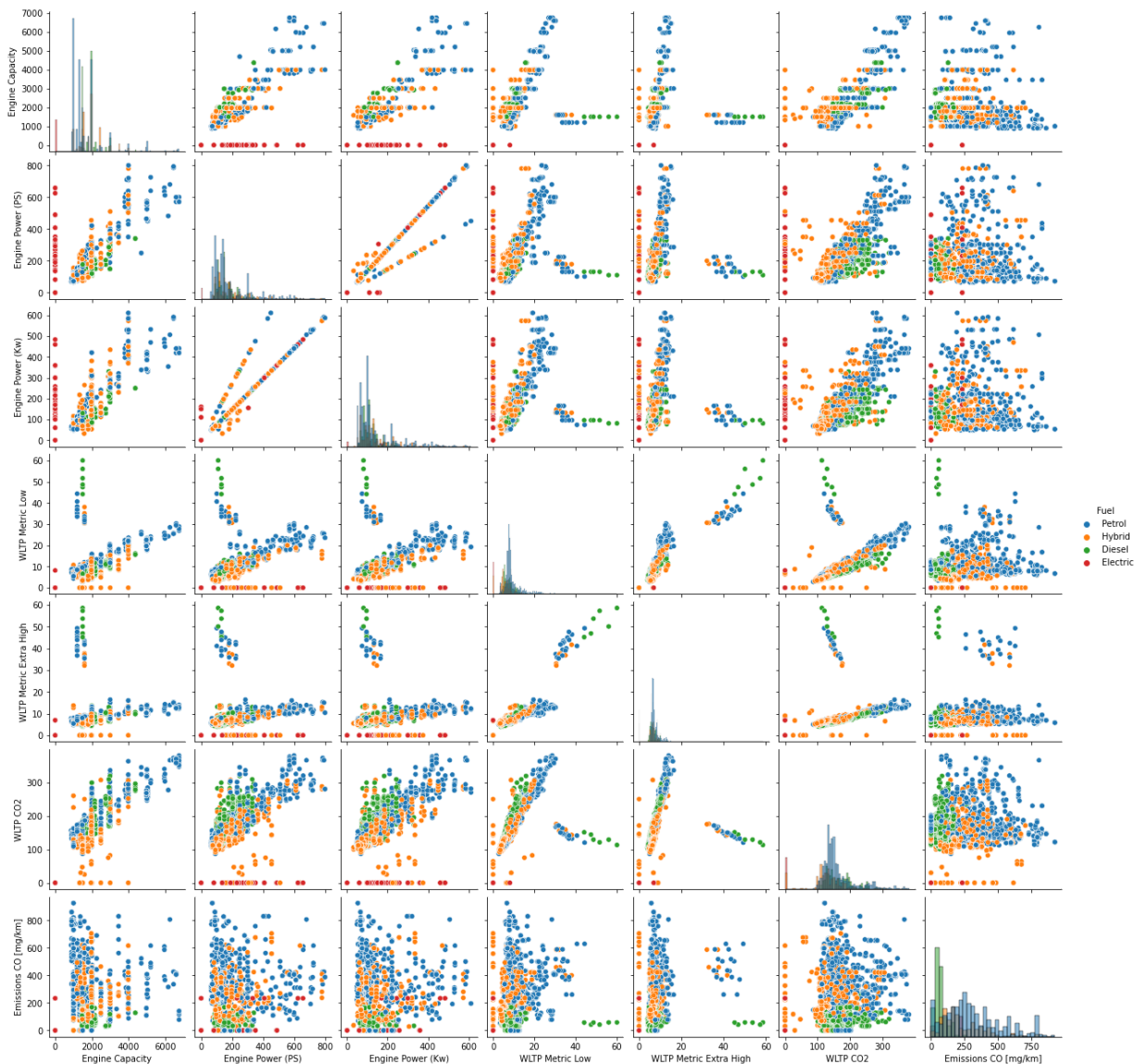
Next, we look at the pairplot of some of the main features. It is clear that hybrid and electric cars act like outliers which is because of their low emissions and engine size related figures. Generally speaking, the dataset has many outliers and there is no obvious correlation among the features and CO Emission

It is also clear petrol and diesel have higher WLTP figures than the rest.

```
In [24]: sns.pairplot(df,vars=['Engine Capacity','Engine Power (PS)','Engine Power (Kw',
                             'WLTP Metric Low','WLTP Metric Extra High','WLTP CO2',
```

```
'Emissions CO [mg/km]'], hue='Fuel',
diag_kind='hist')
```

Out[24]: <seaborn.axisgrid.PairGrid at 0x7fec5d00e340>



4.5 - Mean and STD

Note how each feature covers a very different range, therefore we need to normalise the dataset when building the prediction models.

```
In [25]: df.describe().transpose()[['mean', 'std']]
```

Out[25]:

	mean	std
Engine Capacity	1744.287665	894.872988
Engine Power (PS)	175.893217	109.439872
Engine Power (Kw)	131.782628	83.890608
WLTP Metric Low	8.307146	4.909390
WLTP Metric Medium	6.399704	4.837777
WLTP Metric High	5.861178	4.419764
WLTP Metric Extra High	7.056735	3.763594

	mean	std
WLTP Metric Combined	6.737001	4.058259
WLTP CO2	151.679009	57.911991
Emissions CO [mg/km]	232.946411	197.198821

5 Encoding

Let's look at the result df and see how many categorical features we have:

```
In [26]: for column in df:
          print("{} | {} | {}".format(
              df[column].name, len(df[column].unique()), df[column].dtype
          ))
```

```
Manufacturer | 40 | object
Model | 305 | object
Description | 2138 | object
Engine Capacity | 88 | float64
Engine Power (PS) | 196 | float64
Engine Power (Kw) | 179 | float64
WLTP Metric Low | 240 | float64
WLTP Metric Medium | 168 | float64
WLTP Metric High | 132 | float64
WLTP Metric Extra High | 138 | float64
WLTP Metric Combined | 146 | float64
WLTP CO2 | 246 | float64
Emissions CO [mg/km] | 422 | float64
Transmission | 3 | object
Fuel | 4 | object
Powertrain | 3 | object
```

There are three categorical features that we are going to include in our analysis. These are:

- Transmission
- Fuel
- Powertrain

We are using `get_dummies()` for encoding those features.

We also going to drop Manufacturer, Model, and Description. They are not useful for prediction.

```
In [27]: dummy_cols = ['Transmission', 'Fuel', 'Powertrain']
df_encode = pd.get_dummies(df, columns=dummy_cols).copy()
df_encode = df_encode.drop(['Manufacturer', 'Model', 'Description'], axis=1).copy()
```

And finally we check data types do we have and dim of our final dataset.

```
In [28]: print(df_encode.dtypes)
          print()
          print('Dataset dim is ', df_encode.shape)
```

```
Engine Capacity      float64
Engine Power (PS)    float64
Engine Power (Kw)    float64
WLTP Metric Low      float64
WLTP Metric Medium   float64
```

WLTP Metric High	float64
WLTP Metric Extra High	float64
WLTP Metric Combined	float64
WLTP CO2	float64
Emissions CO [mg/km]	float64
Transmission_Automatic	uint8
Transmission_Electric	uint8
Transmission_Manual	uint8
Fuel_Diesel	uint8
Fuel_Electric	uint8
Fuel_Hybrid	uint8
Fuel_Petrol	uint8
Powertrain_EV	uint8
Powertrain_Hybrid	uint8
Powertrain_ICE	uint8
dtype:	object

Dataset dim is (4467, 20)

6 Prediction Models

6.1 - SVM Model - Binary Classification

We build a binary classification model to predict if the car produce high or low CO emission.
</br>

The new binary target classes are 0 or 1. Any sample data above the mean of the original target will belong to class 0 and considered as high (bad). On the other hand, those below its mean are 1 and considered as low (good).

```
In [29]: # Mean of the target value.
m = df_encode['Emissions CO [mg/km]'].mean()

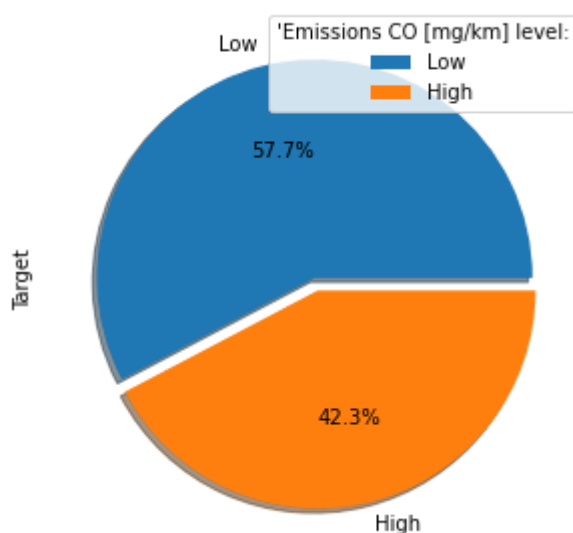
# make a copy of df_encode
df_svm = df_encode.copy()
df_svm['Target'] = (df_svm['Emissions CO [mg/km]'] <= m).astype('int64')
```

Plot Percentages Of Classes

```
In [30]: df_svm['Target'].value_counts() / df_svm['Target'].count()
y_labels = ["Low", "High"]
df_svm['Target'].value_counts().plot.pie(explode=[0, 0.06],
                                         figsize=(5,5),
                                         labels = y_labels,
                                         autopct = '%1.1f%%', shadow=True)

# The pie plot shows proportion of each class of the target value.
plt.legend(title = "'Emissions CO [mg/km] level:")
```

Out[30]: <matplotlib.legend.Legend at 0x7fec5fbb12e0>



Define Inputs And Target Variables

Our target variable is the new feature Target which contains 2 classes 0 and 1. All other numeric features are our input variables.

```
In [31]: y_2 = df_svm['Target']
X_2 = df_svm.drop(['Emissions CO [mg/km]', 'Target'], axis=1).copy()
print('Our SVM Input Features\n')
X_2.dtypes
```

Our SVM Input Features

```
Out[31]: Engine Capacity          float64
Engine Power (PS)          float64
Engine Power (Kw)          float64
WLTP Metric Low            float64
WLTP Metric Medium         float64
WLTP Metric High           float64
WLTP Metric Extra High     float64
WLTP Metric Combined       float64
WLTP CO2                   float64
Transmission_Automatic     uint8
Transmission_Electric      uint8
Transmission_Manual        uint8
Fuel_Diesel                uint8
Fuel_Electric              uint8
Fuel_Hybrid                uint8
Fuel_Petrol                uint8
Powertrain_EV              uint8
Powertrain_Hybrid          uint8
Powertrain_ICE             uint8
dtype: object
```

Splitting The Dataset Into Train And Test

Split train set and test set into 80% to 20%.

```
In [32]: X_train_2, X_test_2, y_train_2, y_test_2 = train_test_split(X_2, y_2, test_si
```

Build SVM Model

We are using MinMaxScaler for normalisation of both train and test data.</br>

```
In [33]: from sklearn.preprocessing import MinMaxScaler
from sklearn import svm
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import plot_confusion_matrix
import sklearn.metrics as metrics

#MinMax normalization
scaler = MinMaxScaler()

# Scaling the data
X_train_2 = scaler.fit_transform(X_train_2)
X_test_2 = scaler.transform(X_test_2)
```

Check feature scalling

```
In [34]: df_describe_2 = pd.DataFrame(X_train_2)
df_describe_2.describe().transpose()
```

```
Out[34]:
```

	count	mean	std	min	25%	50%	75%	max
0	3573.0	0.258063	0.131438	0.0	0.197066	0.236776	0.295896	1.0
1	3573.0	0.220430	0.136284	0.0	0.142500	0.181250	0.250000	1.0
2	3573.0	0.216014	0.136122	0.0	0.137255	0.174837	0.240196	1.0
3	3573.0	0.148184	0.087956	0.0	0.106952	0.135472	0.163993	1.0
4	3573.0	0.101402	0.069361	0.0	0.077901	0.092210	0.111288	1.0
5	3573.0	0.082307	0.063586	0.0	0.064426	0.075630	0.088235	1.0
6	3573.0	0.123290	0.067509	0.0	0.102967	0.116928	0.137871	1.0
7	3573.0	0.110240	0.067885	0.0	0.086601	0.102941	0.119281	1.0
8	3573.0	0.402077	0.154681	0.0	0.336870	0.384615	0.464191	1.0
9	3573.0	0.620767	0.485264	0.0	0.000000	1.000000	1.000000	1.0
10	3573.0	0.006437	0.079985	0.0	0.000000	0.000000	0.000000	1.0
11	3573.0	0.372796	0.483616	0.0	0.000000	0.000000	1.000000	1.0
12	3573.0	0.238735	0.426370	0.0	0.000000	0.000000	0.000000	1.0
13	3573.0	0.034985	0.183767	0.0	0.000000	0.000000	0.000000	1.0
14	3573.0	0.213826	0.410063	0.0	0.000000	0.000000	0.000000	1.0
15	3573.0	0.512455	0.499915	0.0	0.000000	1.000000	1.000000	1.0
16	3573.0	0.034985	0.183767	0.0	0.000000	0.000000	0.000000	1.0
17	3573.0	0.240134	0.427225	0.0	0.000000	0.000000	0.000000	1.0
18	3573.0	0.724881	0.446637	0.0	0.000000	1.000000	1.000000	1.0

Now we have the data ready for the model.

Create a SVM classifier and use linear as kernel as it is the most common one. First we make prediction with train data, and we also going to look at the accuracy of the model by printing the metrics report and the confusion matrix.

```
In [35]: svm_1 = svm.SVC(kernel='linear',C=1)

# Fit the model
svm_1.fit(X_train_2, y_train_2)

# Make predictions on train dataset
y_pred_2 = svm_1.predict(X_train_2)
```

Model Evaluation

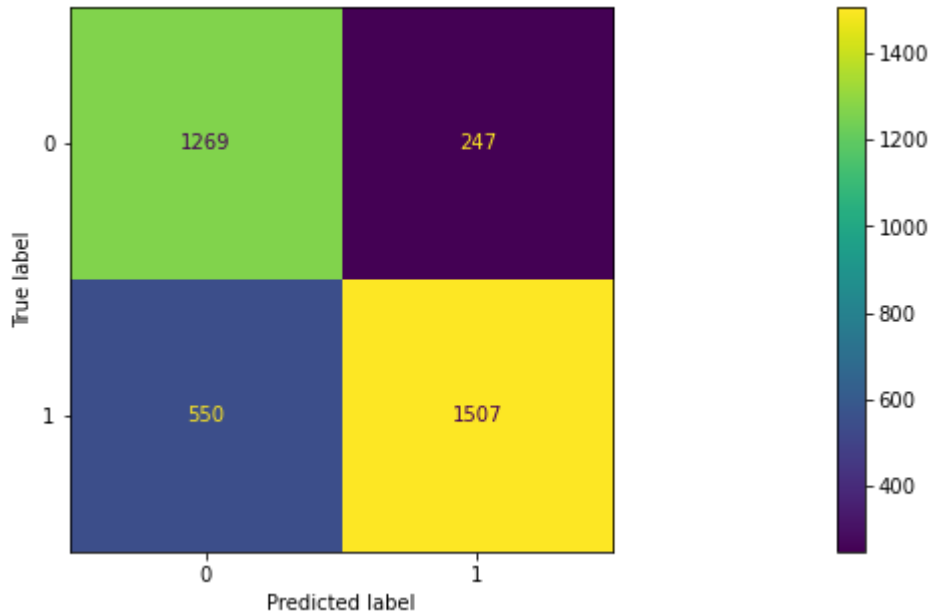
10-fold cross validation is used here to check the accuracy of the model.

```
In [36]: accuracies = cross_val_score(svm_1, X_train_2, y_train_2, cv=10)
print("Train Score:", np.mean(accuracies))
print("confusion_matrix:", confusion_matrix(y_train_2, y_pred_2))
print(metrics.classification_report(y_train_2, y_pred_2, digits=2))
plot_confusion_matrix(svm_1, X_train_2, y_train_2)
plt.show()
```

Train Score: 0.774972223526282

confusion_matrix: [[1269 247]
[550 1507]]

	precision	recall	f1-score	support
0	0.70	0.84	0.76	1516
1	0.86	0.73	0.79	2057
accuracy			0.78	3573
macro avg	0.78	0.78	0.78	3573
weighted avg	0.79	0.78	0.78	3573



Average score for training data is 77%.

Precision figure for class 0 is 70% and for class 1 is 86% which is not bad. Out of 1516 sample with class 0, 1269 are correct and out of 2057 of class 1, 1507 are correct.

Let's test accuracy with test dataset.

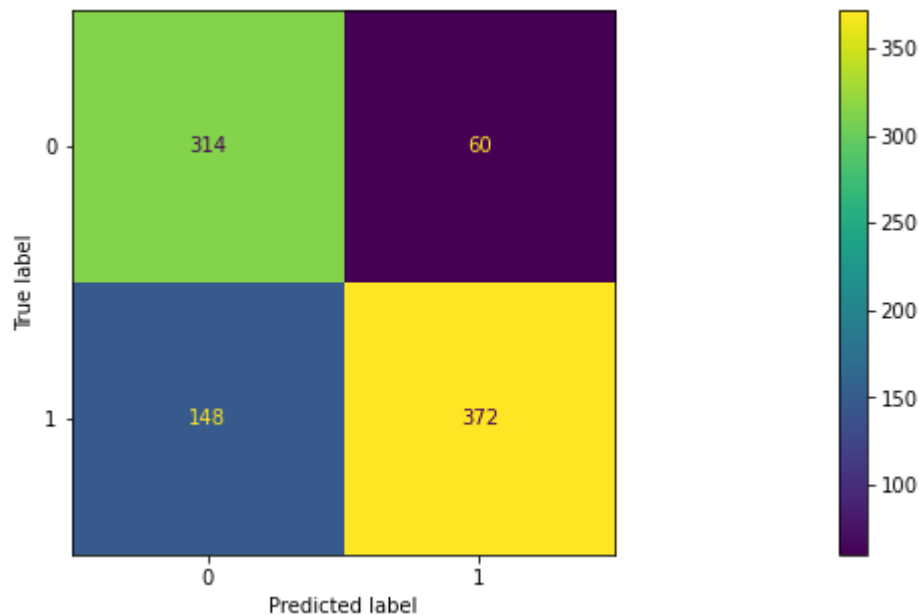
```
In [37]: y_pred_3 = svm_1.predict(X_test_2)

#Accuracies of the model. 10-fold cross validation is used here.
accuracies2 = cross_val_score(svm_1, X_test_2, y_test_2, cv=10)
print("Test Score:", np.mean(accuracies2))
print("confusion_matrix:", confusion_matrix(y_test_2, y_pred_3))
print(metrics.classification_report(y_test_2, y_pred_3, digits=2))
plot_confusion_matrix(svm_1, X_test_2, y_test_2)
plt.show()
```

Test Score: 0.7393008739076155

confusion_matrix: [[314 60]
[148 372]]

	precision	recall	f1-score	support
0	0.68	0.84	0.75	374
1	0.86	0.72	0.78	520
accuracy			0.77	894
macro avg	0.77	0.78	0.77	894
weighted avg	0.79	0.77	0.77	894



Average score of test data is 74% which is less than training data. It is normal to have less test score than train score.

precision figures are close with 68% for class 0 and 86% for class 1.

The score is not that great and that could be because there are many outliers in our dataset as you can see from the pairplot figure above.

We are going to use GridSearchCV to find the best hyperparameters C, kernel, degree, and gamma.

GridSearchCV Model

GridSearchCV is used here to search the best parameters for the SVM model.

In [38]:

```
grid = {
    'C':[0.01, 0.1, 1, 10],
    'kernel' : ["linear", "poly", "rbf", "sigmoid"],
    'degree' : [1, 3, 5, 7],
    'gamma' : [0.01, 1]
}

svm_2 = SVC(random_state = 125)
svm_cv = GridSearchCV(svm_2, grid, cv = 10)
svm_cv.fit(X_train_2, y_train_2)

print("Best Parameters:", svm_cv.best_params_)
print("Accuracy on train data", svm_cv.best_score_)
# Print the accuracy on the test data
print("Accuracy on test data:", svm_cv.score(X_test_2, y_test_2))
```

Best Parameters: {'C': 10, 'degree': 7, 'gamma': 1, 'kernel': 'poly'}
Accuracy on train data 0.8804868316041501
Accuracy on test data: 0.8780760626398211

We use the best parameters for prediction and print the classification report.

In [39]:

```
grid_predictions = svm_cv.predict(X_test_2)

print(classification_report(y_test_2, grid_predictions))
```

	precision	recall	f1-score	support
0	0.82	0.90	0.86	374
1	0.92	0.86	0.89	520
accuracy			0.88	894
macro avg	0.87	0.88	0.88	894
weighted avg	0.88	0.88	0.88	894

The precision figures for both classes are better than the linear SVM model. Class 0 was 68% now it is 82%. Class 1 was 86% and now it is 92%.

Use best parameters for prediction

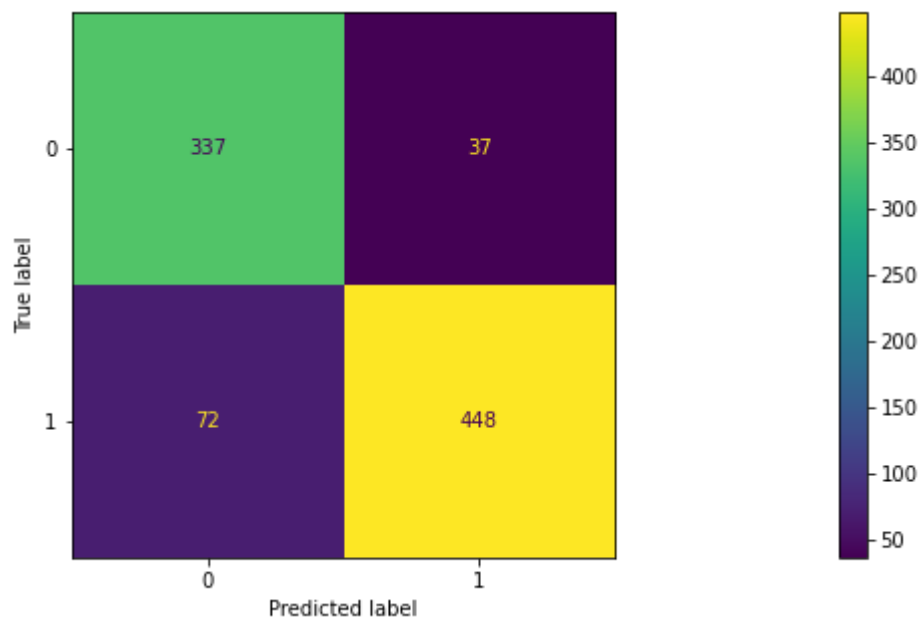
In [40]:

```
svm_3 = SVC(kernel='poly', degree=7, gamma=1, C=10)
svm_3.fit(X_train_2, y_train_2)
y_pred_2_1 = svm_3.predict(X_test_2)

print("Test Accuracy:\n", accuracy_score(y_test_2, y_pred_2_1))
print("confusion_matrix:\n", confusion_matrix(y_test_2, y_pred_2_1))
print(metrics.classification_report(y_test_2, y_pred_2_1, digits=2))
plot_confusion_matrix(svm_3, X_test_2, y_test_2)
plt.show()
```

```
Test Accuracy:
0.8780760626398211
confusion_matrix:
[[337  37]
 [ 72 448]]
```

	precision	recall	f1-score	support
0	0.82	0.90	0.86	374
1	0.92	0.86	0.89	520
accuracy			0.88	894
macro avg	0.87	0.88	0.88	894
weighted avg	0.88	0.88	0.88	894



Accuracy score for the new polynomial model has also improved. Out of 374 data samples with class 0, only 37 were misclassified. And out of 520 sample with class 1, only 72 were wrongly classified as 0.

Overall, this binary SVM model is good and produces good results.

6.2 - Neural Network Regression Model

Build the model using Keras deep learning API <https://keras.io/>

Identify The Input And Target Variables

Our target variable (y) is 'Emissions CO [mg/km]' which is continuous variable. So our NN model is regression.

We use all the numeric features for the prediction models as input variables or X.

```
In [41]: X = df_encode.loc[:, df_encode.columns != 'Emissions CO [mg/km]']
print()
print('X or input variables are:')
print()
print(X.columns.tolist())
y = df_encode['Emissions CO [mg/km]']

X or input variables are:

['Engine Capacity', 'Engine Power (PS)', 'Engine Power (Kw)', 'WLTP Metric Low', 'WLTP Metric Medium', 'WLTP Metric High', 'WLTP Metric Extra High', 'WLTP Metric Combined', 'WLTP CO2', 'Transmission_Automatic', 'Transmission_Electric', 'Transmission_Manual', 'Fuel_Diesel', 'Fuel_Electric', 'Fuel_Hybrid', 'Fuel_Petrol', 'Powertrain_EV', 'Powertrain_Hybrid', 'Powertrain_ICE']
```

Splitting The Data

```
In [42]: X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=.2)
```

Define Keras Sequential Model

The model is using build-in Adam() optimiser.

For measuring the losses, we use mean_square_error metrics.

We are using sklearn.preprocessing.StandardScaler library to scale our train and test data.

Our model has seven hidden layers and one output layer with linear activation function.

```
In [43]: from keras.models import Sequential
#from keras import utils
import keras
from keras.layers import Dense, Activation
import scikeras
from scikeras.wrappers import KerasRegressor
from sklearn.model_selection import cross_val_score
from keras.optimizers import SGD
from keras import initializers
from keras.layers import LeakyReLU

numerics = ['uint8', 'int16', 'int32', 'int64', 'float16', 'float32', 'float64']

X_train_chris = X_train.select_dtypes(include=numerics).to_numpy()
X_test_chris = X_test.select_dtypes(include=numerics).to_numpy()
```

```

y_train_chris = Y_train.to_numpy()
y_test_chris = Y_test.to_numpy()

#Standard normalization
sc = StandardScaler()
X_train_chris = sc.fit_transform(X_train_chris)
X_test_chris = sc.transform(X_test_chris)
y_train_chris = sc.fit_transform(y_train_chris.reshape(len(y_train_chris),1))
y_test_chris = sc.transform(y_test_chris.reshape(len(y_test_chris),1))[:,0]

dim = X_train_chris.shape[1]

# define the keras model
model = Sequential()

# First layer with inputs
model.add(Dense(64, input_shape=(dim,), activation='relu'))

# Second hidden layer
model.add(Dense(32, activation='relu'))

# Third hidden layer
model.add(Dense(132, activation='relu'))

# Fourth hidden layer
model.add(Dense(32))
model.add(LeakyReLU(alpha=0.1))

# Fifth hidden layer
model.add(Dense(32))
model.add(LeakyReLU(alpha=0.1))

# Sixth hidden layer
model.add(Dense(16))
model.add(LeakyReLU(alpha=0.1))

# Seventh hidden layer
model.add(Dense(4))
model.add(LeakyReLU(alpha=0.1))

# Output layer is linear
model.add(Dense(1, activation='linear'))

opt = keras.optimizers.Adam(learning_rate=0.001)
# compile the keras model with adam optimise

model.compile(loss='mean_squared_error', optimizer=opt , metrics=['mean_squared_error'])
model.summary()

```

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 64)	1280
dense_1 (Dense)	(None, 32)	2080
dense_2 (Dense)	(None, 132)	4356
dense_3 (Dense)	(None, 32)	4256
leaky_re_lu (LeakyReLU)	(None, 32)	0
dense_4 (Dense)	(None, 32)	1056

leaky_re_lu_1 (LeakyReLU)	(None, 32)	0
dense_5 (Dense)	(None, 16)	528
leaky_re_lu_2 (LeakyReLU)	(None, 16)	0
dense_6 (Dense)	(None, 4)	68
leaky_re_lu_3 (LeakyReLU)	(None, 4)	0
dense_7 (Dense)	(None, 1)	5

=====

Total params: 13,629
 Trainable params: 13,629
 Non-trainable params: 0

Fit The Model

```
In [44]: history = model.fit(X_train_chris, y_train_chris,
                             validation_data=(X_test_chris, y_test_chris), epochs=1200, verbose=0)
```

Evaluate The Model

We are using MSE as metric to measure the accuracy of our NN model.

```
In [45]: train_mse, train_accuracy = model.evaluate(X_train_chris, y_train_chris, verbose=0)
test_mse, test_accuracy = model.evaluate(X_test_chris, y_test_chris, verbose=0)
print('Train MSE: %.5f, Test MSE: %.5f' % (train_mse, test_mse))
```

Train MSE: 0.01274, Test MSE: 0.13673

Plot Train And Test Losses

```
In [46]: plt.figure(figsize=(14,4))
plt.title('Loss / Mean Squared Error')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.ylabel('MSE')
plt.xlabel('Epochs')
plt.yticks(np.arange(0, np.max(history.history['loss'])+0.2, step=0.2))
plt.legend()
plt.show()
```



Keras Model Evaluation

The Loss curves for both training and test data are fluctuating and not smooth. This could be because we have relatively small dataset and little large network (13,629 parameters).

Our train MSE is on average 1% while the test MSE is 8%-17%. It is normal to training MSE lower than test MSE.

This model prediction accuracy is very high and therefore it is a good model.

6.3 - Lasso Regression Model

First, we import the libraries we are going to use to build the Lasso regression model.

```
In [47]: from sklearn.linear_model import Lasso, LassoCV
from sklearn.preprocessing import scale
from sklearn.model_selection import RepeatedKFold
from sklearn.metrics import mean_squared_error
```

Creating the Training and Test Datasets

```
In [48]: X_train_lasso, X_test_lasso, y_train_lasso, y_test_lasso = train_test_split(X
```

Data Normalisation

We use StandardScaler() function for scalling the datasets.

```
In [49]: sc = StandardScaler()
y_train_lasso = y_train_lasso.to_numpy()
y_test_lasso = y_test_lasso.to_numpy()
X_train_lasso = sc.fit_transform(X_train_lasso)
X_test_lasso = sc.transform(X_test_lasso)
y_train_lasso = sc.fit_transform(y_train_lasso.reshape(len(y_train_lasso),1))
y_test_lasso = sc.transform(y_test_lasso.reshape(len(y_test_lasso),1))[:,0]
```

Exploring L1 Penalty Values

First, we investigate the size of each feature weight as function of alpha.

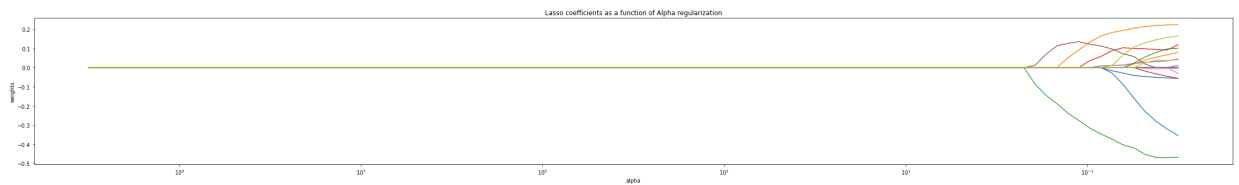
```
In [50]: alphas = 10**np.linspace(10,-2,100)*0.5

# Define the model
lasso = Lasso()
coefs = []

for a in alphas*2:
    lasso.set_params(alpha=a)
    lasso.fit(X_train_lasso, y_train_lasso)
    coefs.append(lasso.coef_)

ax = plt.gca()
ax.plot(alphas*2, coefs)
ax.set_xscale('log')
```

```
ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
plt.axis('tight')
plt.xlabel('alpha')
plt.ylabel('weights')
plt.title('Lasso coefficients as a function of Alpha regularization');
```



Moving from the left to right in the plot, at first, the coefficient estimates approximate towards zero. Then the model starts to have more predictors with high magnitudes of coefficient estimates.

Selecting Optimal Alpha Value

Next, we need to find the optimal value of alpha to use in our model.

use LassoCV function to fit the regression model and find the optimal alpha. We also use RepeatedKFold() to evaluate the lasso model.

We will define a range for alpha from 0 to 1 with increment of 0.01

```
In [51]: # RepeatedKFold for evaluation of the model
cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1)

# Define the model
lassocv = LassoCV(alphas=np.arange(0.01, 1, 0.01), cv=cv, n_jobs=-1)

# Fit the model
lassocv.fit(X_train_lasso, y_train_lasso)

# Best alpha
best_alpha = lasso cv.alpha_
print('Best lambda that produced the lowest test MSE = ', best_alpha)
print()
print('Coefficients of the model are \n')
pd.Series(lassocv.coef_, index=X.columns)
```

Best lambda that produced the lowest test MSE = 0.01

Coefficients of the model are

```
Out[51]: Engine Capacity          -0.354721
Engine Power (PS)             0.078438
Engine Power (Kw)             0.100952
WLTP Metric Low               0.118721
WLTP Metric Medium            0.045414
WLTP Metric High              -0.000000
WLTP Metric Extra High        -0.030387
WLTP Metric Combined          -0.000000
WLTP CO2                      0.164065
Transmission_Automatic        -0.000000
Transmission_Electric         -0.056444
Transmission_Manual           0.224048
Fuel_Diesel                   -0.468157
Fuel_Electric                 -0.056148
Fuel_Hybrid                   0.011740
Fuel_Petrol                   0.000000
```

```
Powertrain_EV          -0.003684
Powertrain_Hybrid      -0.000000
Powertrain_ICE          0.041792
dtype: float64
```

We can see few of the coefficients are zero value.

Build The Model

Now we are going to use the best lambda in our model for prediction

```
In [52]: lasso.set_params(alpha=lassocv.alpha_)
lasso.fit(X_train_lasso, y_train_lasso)
y_pred = lasso.predict(X_train_lasso)
```

Model Evaluation

```
In [53]: MSE = mean_squared_error(y_train_lasso, y_pred)
print("Train MSE = ", round(MSE, 2) )

y_pred_t = lasso.predict(X_test_lasso)
MSE_test = mean_squared_error(y_test_lasso, y_pred_t)
print("Test MSE = ", round(MSE_test, 2) )
```

```
Train MSE =  0.63
Test MSE =  0.63
```

The model scores 0.63 on the training dataset and 0.67 on the test dataset. The difference is very small but they indicate the model is not really good one.

6.4 - Naive Bayes Model

Additional simple statistical analysis

```
In [54]: df.corr()
```

```
Out[54]:
```

	Engine Capacity	Engine Power (PS)	Engine Power (Kw)	WLTP Metric Low	WLTP Metric Medium	WLTP Metric High	WLTP Metric Extra High	WLTP Metric Combined
Engine Capacity	1.000000	0.774230	0.757315	0.646861	0.399634	0.315369	0.379813	0.413805
Engine Power (PS)	0.774230	1.000000	0.960986	0.575433	0.339775	0.251194	0.272663	0.335780
Engine Power (Kw)	0.757315	0.960986	1.000000	0.578785	0.335136	0.243554	0.257499	0.326848
WLTP Metric Low	0.646861	0.575433	0.578785	1.000000	0.825314	0.829728	0.850635	0.896541
WLTP Metric Medium	0.399634	0.339775	0.335136	0.825314	1.000000	0.878931	0.865129	0.888418

	Engine Capacity	Engine Power (PS)	Engine Power (Kw)	WLTP Metric Low	WLTP Metric Medium	WLTP Metric High	WLTP Metric Extra High	WLTP Metric Combined
WLTP Metric High	0.315369	0.251194	0.243554	0.829728	0.878931	1.000000	0.974594	0.984246
WLTP Metric Extra High	0.379813	0.272663	0.257499	0.850635	0.865129	0.974594	1.000000	0.981181
WLTP Metric Combined	0.413805	0.335780	0.326848	0.896543	0.888418	0.984246	0.981181	1.000000
WLTP CO2	0.742957	0.565200	0.556512	0.775881	0.511176	0.441096	0.557144	0.563379
Emissions CO [mg/km]	-0.109571	0.029086	0.030707	0.156310	0.126220	0.112127	0.118176	0.127399

By checking the correlation, it shows that emissions CO does not have high correlation into other data. This fulfill the Naive Bayes assumption that the attributes are independent from each other. Now, we are going to check the emissions CO data's value because we would like to it from continuous into discrete number by classified it.

```
In [55]: df.describe().transpose()[['mean', 'std']]
```

```
Out[55]:
```

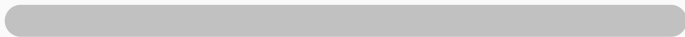
	mean	std
Engine Capacity	1744.287665	894.872988
Engine Power (PS)	175.893217	109.439872
Engine Power (Kw)	131.782628	83.890608
WLTP Metric Low	8.307146	4.909390
WLTP Metric Medium	6.399704	4.837777
WLTP Metric High	5.861178	4.419764
WLTP Metric Extra High	7.056735	3.763594
WLTP Metric Combined	6.737001	4.058259
WLTP CO2	151.679009	57.911991
Emissions CO [mg/km]	232.946411	197.198821

```
In [56]: df[df['Emissions CO [mg/km]']==df['Emissions CO [mg/km]'].max()]
```

```
Out[56]:
```

	Manufacturer	Model	Description	Engine Capacity	Engine Power (PS)	Engine Power (Kw)	WLTP Metric Low	WLTP Metric Medium	WLTP Metric High
3332	RENAULT	Clio	Iconic SCe 75	999.0	72.0	53.0	6.7	5.0	4.7

	Manufacturer	Model	Description	Engine Capacity	Engine Power (PS)	Engine Power (Kw)	WLTP Metric Low	WLTP Metric Medium	WLTP Metric High
3333	RENAULT	Clio	Iconic SCe 75	999.0	72.0	53.0	6.7	5.1	4.8
3334	RENAULT	Clio	Iconic SCe 75 with BOSE	999.0	72.0	53.0	6.7	5.0	4.7
3335	RENAULT	Clio	Iconic SCe 75 with BOSE	999.0	72.0	53.0	6.7	5.1	4.8
3382	RENAULT	Clio	Play SCe 75	999.0	72.0	53.0	6.7	5.0	4.7
3383	RENAULT	Clio	Play SCe 75	999.0	72.0	53.0	6.7	5.1	4.7



In [57]:

```
df[df['Emissions CO [mg/km]']==df['Emissions CO [mg/km]'].min()]
```

Out[57]:

	Manufacturer	Model	Description	Engine Capacity	Engine Power (PS)	Engine Power (Kw)	WLTP Metric Low	WLTP Metric Medium	WLTP Metric High
256	CITROEN	New C4	100kW Electric Vehicle with 50kWh battery	0.0	136.0	100.0	0.0	0.0	0.0
257	CITROEN	New C4	100kW Electric Vehicle with 50kWh battery	0.0	136.0	100.0	0.0	0.0	0.0
270	CITROEN	SpaceTourer	50KWh Electric Vehicle	0.0	136.0	100.0	0.0	0.0	0.0
271	CITROEN	SpaceTourer	50KWh Electric Vehicle	0.0	136.0	100.0	0.0	0.0	0.0
704	DS	DS 3 CROSSBACK	E-TENSE	0.0	0.0	0.0	0.0	0.0	0.0
...
4245	TESLA	Model 3	Performance (E5D#Gp)	0.0	0.0	0.0	0.0	0.0	0.0
4246	TESLA	Model S	Long Range (SA3EB)	0.0	0.0	0.0	0.0	0.0	0.0
4247	TESLA	Model S	Performance (SA3EP)	0.0	0.0	0.0	0.0	0.0	0.0
4248	TESLA	Model X	Long Range (XA3EB)	0.0	0.0	0.0	0.0	0.0	0.0
4249	TESLA	Model X	Performance (XA3EP)	0.0	0.0	0.0	0.0	0.0	0.0

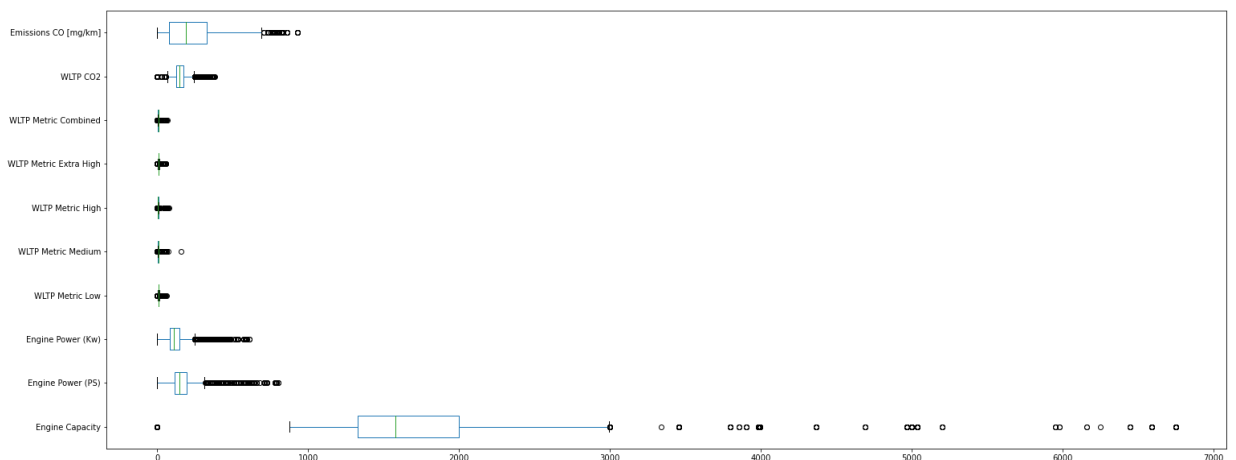
283 rows × 16 columns

```
In [58]: df['Emissions CO [mg/km]'].value_counts()
```

```
Out[58]: 0.000000    283
233.791603     93
96.000000     91
71.000000     76
88.000000     73
...
368.000000     1
131.000000     1
255.000000     1
808.000000     1
38.000000      1
Name: Emissions CO [mg/km], Length: 422, dtype: int64
```

```
In [59]: plt.rcParams["figure.figsize"] = (25,10)
df.plot.box(vert=False)
```

```
Out[59]: <AxesSubplot:>
```



The analysis above are several consideration to split the data. The maximum emissions CO is 927 mg/km while minimum is 0. The data distribution is not well spread as the median is more in the left, leaving several threshold on the left. The mean of the data is 232.946411.

Converting the Emissions CO data into Discrete Number

Here we will classify the amount of emissions CO into zero to low emissions and high emissions.

```
In [60]: df['Emissions CO Class'] = pd.qcut(df['Emissions CO [mg/km]'], 2, labels=['zero to low', 'high'])
```

```
In [61]: df['Emissions CO Class'].value_counts()
```

```
Out[61]: zero to low    2237
high                2230
Name: Emissions CO Class, dtype: int64
```

The data is well separated with almost same number. Next, we will start to do the Naive Bayes methods.

Naive bayes method

This method is originally comes from Bayes' Theorem. It is naive, because it assumes that all of the attributes are independent from each other. By calculating the conditional probabilities, there will be a likelihood probability, prior probability and evidence as considerations to predict the probability of something happens. The equation is shown as below:

$$P(C_k|x) = \frac{P(C_k) * P(x|C_k)}{P(x)}$$

where $P(C_k|x)$ is the posterior probability, followed by likelihood and prior probability as numerator and evidence probability as denominator.

In our case, we will determine how well the method to predict each class of Emissions CO. Because the data is in continuous number, we will convert it as discrete number. We classify the data into 'zero to low'and 'high' emissions. It is followed by cross-validation to avoid over-fitting.

```
In [62]: from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OrdinalEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.naive_bayes import CategoricalNB
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.model_selection import cross_val_score
from sklearn.metrics import roc_curve, auc
```

```
In [63]: df['High CO Emission']=df['Emissions CO Class'].apply(lambda x: 1 if x=='high
```

```
In [64]: df['Emissions CO Class'].value_counts()
```

```
Out[64]: zero to low    2237
high              2230
Name: Emissions CO Class, dtype: int64
```

```
In [65]: df.iloc[:,12:]
```

```
Out[65]:
```

	Emissions CO [mg/km]	Transmission	Fuel	Powertrain	Emissions CO Class	High CO Emission
0	760.0	Manual	Petrol	ICE	high	1
1	760.0	Manual	Petrol	ICE	high	1
2	760.0	Manual	Petrol	ICE	high	1
3	760.0	Manual	Petrol	ICE	high	1
4	829.0	Manual	Petrol	ICE	high	1
...
4652	423.0	Automatic	Hybrid	Hybrid	high	1
4653	423.0	Automatic	Hybrid	Hybrid	high	1

	Emissions CO [mg/km]	Transmission	Fuel	Powertrain	Emissions CO Class	High CO Emission
4654	706.0	Automatic	Hybrid	Hybrid	high	1
4655	706.0	Automatic	Hybrid	Hybrid	high	1
4656	706.0	Automatic	Hybrid	Hybrid	high	1

4467 rows × 6 columns

Gaussian Naive Bayes

It assumes that every countinuous values are distributed into a normal(Gaussian) distribution.

```
In [66]: df['Fuel'].value_counts()
```

```
Out[66]: Petrol      2307
         Diesel      1068
         Hybrid       948
         Electric     144
         Name: Fuel, dtype: int64
```

```
In [67]: df['Fuel Class']=df['Fuel'].apply(lambda x: 1 if x=='Petrol' else
                                           0 if x=='Diesel' or 'Hybrid' else -1)
```

```
In [68]: x = df[['Fuel Class', 'WLTP CO2']]
         y = df['High CO Emission']

         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [69]: model = GaussianNB()
         clf = model.fit(X_train, y_train)

         pred_labels = model.predict(X_test)

         print(classification_report(y_test, pred_labels))
```

	precision	recall	f1-score	support
0	0.75	0.68	0.71	455
1	0.70	0.76	0.73	439
accuracy			0.72	894
macro avg	0.72	0.72	0.72	894
weighted avg	0.72	0.72	0.72	894

```
In [70]: model.score(X_test, y_test)
```

```
Out[70]: 0.7192393736017897
```

```
In [71]: cross_val_score(GaussianNB(), X, y, cv = 5)
```

```
Out[71]: array([0.8557047 , 0.69686801, 0.64277716, 0.71444569, 0.68085106])
```

Categorical Naive Bayes

```
In [72]: x = df[['Transmission', 'Powertrain']]
y = df['High CO Emission'].values

enc = OrdinalEncoder()
X = enc.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = CategoricalNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
```

	precision	recall	f1-score	support
0	0.59	0.72	0.65	455
1	0.62	0.48	0.54	439
accuracy			0.60	894
macro avg	0.61	0.60	0.59	894
weighted avg	0.61	0.60	0.60	894

```
In [73]: print('Accuracy score: ', clf.score(X_test, y_test))
```

Accuracy score: 0.6017897091722595

```
In [74]: cross_val_score(CategoricalNB(), X, y, cv = 5)
```

```
Out[74]: array([0.60514541, 0.74272931, 0.55991041, 0.45128779, 0.61254199])
```

Mixed Naive Bayes

Here we will combine between the continuous attributes and categorical attributes.

```
In [75]: df['WLTPCO2_qt'] = pd.qcut(df['WLTP CO2'], 5, labels=['bottom 20', 'lower 20', 'middle 20', 'upper 20', 'top 20'])

X = df[['Transmission', 'Powertrain', 'Fuel Class', 'WLTPCO2_qt']]
y = df['High CO Emission']

enc = OrdinalEncoder()
X = enc.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = CategoricalNB()
clf = model.fit(X_train, y_train)

label = model.predict(X_test)

print(classification_report(y_test, label))
print('Accuracy score: ', clf.score(X_test, y_test))
```

	precision	recall	f1-score	support
0	0.75	0.68	0.71	455
1	0.70	0.76	0.73	439
accuracy			0.72	894
macro avg	0.72	0.72	0.72	894

weighted avg 0.72 0.72 0.72 894

Accuracy score: 0.7192393736017897

Here we will combine the model fitting of each Gaussian and Categorical. The new probability of each model is combined and to be used into the further fit modelling.

In [76]:

```
XG = df[['Fuel Class', 'WLTP CO2']]
XC = df[['Transmission', 'Powertrain']]

enc = OrdinalEncoder()
XC = enc.fit_transform(XC)

X=np.c_[XG, XC[:,0].ravel(), XC[:,1].ravel()]

Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.2, random_s

# Model Fitting
# Gaussian model
mG = GaussianNB()
fitG = mG.fit(Xtrain[:,0:2], ytrain)
# Categorical model
mC = CategoricalNB()
fitC = mC.fit(Xtrain[:,2:4], ytrain)

# Probability each train and test dataset
# On training data
G_train_probas = mG.predict_proba(Xtrain[:,0:2])
C_train_probas = mC.predict_proba(Xtrain[:,2:4])
# And on testing data
G_test_probas = mG.predict_proba(Xtest[:,0:2])
C_test_probas = mC.predict_proba(Xtest[:,2:4])

# New combination probability
X_new_train = np.c_[(G_train_probas[:,1], C_train_probas[:,1])] # Train
X_new_test = np.c_[(G_test_probas[:,1], C_test_probas[:,1])] # Test

# Fit the model
mG = GaussianNB()
fitG = mG.fit(X_new_train, y_train)

# Predict class labels on a test data
lpred = mG.predict(X_new_test)

print('Classes: ', fitG.classes_) # class labels known to the classifier
print('Class Priors: ', fitG.class_prior_) # probability of each class.
# Use score method to get accuracy of model
print('-----')
score = mG.score(X_new_test, y_test)
print('Accuracy Score: ', score)
print('-----')
# Look at classification report to evaluate the model
print(classification_report(y_test, pred_labels))
```

Classes: [0 1]

Class Priors: [0.49874055 0.50125945]

Accuracy Score: 0.7315436241610739

precision recall f1-score support

	0	0.59	0.72	0.65	455
	1	0.62	0.48	0.54	439
accuracy				0.60	894
macro avg		0.61	0.60	0.59	894
weighted avg		0.61	0.60	0.60	894

From the combination above, the highest accuracy score is obtained in the last method which is combination between Gaussian and Categorical NB. This model is categorized as good model with accuracy higher than 70%, the accuracy is 73.15%.

To improve the accuracy, we try to vary the emission CO into three label as below.

Three Class labels for Emission CO

```
In [77]: df['Emissions CO 3Class'] = pd.qcut(df['Emissions CO [mg/km]'], 3, labels=['low', 'high', 'medium'])
df['Emissions CO 3Class'].value_counts()
```

```
Out[77]: low      1511
high      1481
medium    1475
Name: Emissions CO 3Class, dtype: int64
```

```
In [78]: df['CO Emission'] = df['Emissions CO 3Class'].apply(lambda x: 1 if x=='high' else 0 if x=='medium' else -1)
```

```
In [79]: # Gaussian

X = df[['Fuel Class', 'WLTP CO2']]
y = df['CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)
```

	precision	recall	f1-score	support
-1	0.57	0.70	0.63	305
0	0.53	0.07	0.12	283
1	0.51	0.81	0.63	306
accuracy			0.54	894
macro avg	0.54	0.53	0.46	894
weighted avg	0.54	0.54	0.47	894

Accuracy Score: 0.5391498881431768

```
In [80]: X = df[['Transmission', 'Powertrain']]
y = df['CO Emission'].values

enc = OrdinalEncoder()
X = enc.fit_transform(X)
```



```

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = CategoricalNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)

```

	precision	recall	f1-score	support
-1	0.47	0.59	0.52	305
0	0.52	0.39	0.45	283
1	0.45	0.44	0.45	306
accuracy			0.48	894
macro avg	0.48	0.47	0.47	894
weighted avg	0.48	0.48	0.47	894

Accuracy Score: 0.47651006711409394

By using three classes prediction, the accuracy decrease almost 30%. Therefore, we try to use other attributes to predict Emissions CO.

Improving the Accuracy by trying the other attributes

```
In [81]: df.columns
```

```

Out[81]: Index(['Manufacturer', 'Model', 'Description', 'Engine Capacity',
               'Engine Power (PS)', 'Engine Power (Kw)', 'WLTP Metric Low',
               'WLTP Metric Medium', 'WLTP Metric High', 'WLTP Metric Extra High',
               'WLTP Metric Combined', 'WLTP CO2', 'Emissions CO [mg/km]',
               'Transmission', 'Fuel', 'Powertrain', 'Emissions CO Class',
               'High CO Emission', 'Fuel Class', 'WLTPCO2_qt', 'Emissions CO 3Class',
               'CO Emission'],
              dtype='object')

```

```

In [82]: X = df[['Engine Power (PS)', 'WLTP CO2']]
y = df['High CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)

```

	precision	recall	f1-score	support
0	0.54	0.87	0.67	455
1	0.65	0.24	0.35	439
accuracy			0.56	894
macro avg	0.60	0.56	0.51	894
weighted avg	0.60	0.56	0.51	894

Accuracy Score: 0.5637583892617449

```
In [83]: X = df[['Engine Capacity', 'WLTP CO2']]
y = df['High CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)
```

	precision	recall	f1-score	support
0	0.63	0.87	0.73	455
1	0.78	0.47	0.59	439
accuracy			0.67	894
macro avg	0.70	0.67	0.66	894
weighted avg	0.70	0.67	0.66	894

Accuracy Score: 0.6733780760626398

```
In [84]: X = df[['WLTP Metric Combined', 'WLTP CO2']]
y = df['High CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)
```

	precision	recall	f1-score	support
0	0.54	0.88	0.67	455
1	0.64	0.22	0.32	439
accuracy			0.55	894
macro avg	0.59	0.55	0.50	894
weighted avg	0.59	0.55	0.50	894

Accuracy Score: 0.5548098434004475

```
In [85]: X = df[['Engine Power (Kw)', 'WLTP CO2']]
y = df['High CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)
```

	precision	recall	f1-score	support
0	0.54	0.87	0.67	455
1	0.65	0.24	0.35	439
accuracy			0.56	894
macro avg	0.60	0.56	0.51	894
weighted avg	0.60	0.56	0.52	894

Accuracy Score: 0.5637583892617449

In [86]:

```
X = df[['Engine Power (Kw)', 'Engine Capacity']]
y = df['High CO Emission']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = GaussianNB()
clf = model.fit(X_train, y_train)

pred_labels = model.predict(X_test)

print(classification_report(y_test, pred_labels))
score = model.score(X_test, y_test)
print('Accuracy Score: ', score)
```

	precision	recall	f1-score	support
0	0.56	0.90	0.69	455
1	0.73	0.28	0.41	439
accuracy			0.60	894
macro avg	0.64	0.59	0.55	894
weighted avg	0.64	0.60	0.55	894

Accuracy Score: 0.5950782997762863

The best attributes to be used are still the combination of Fuel, WLTP CO2, Transmission, and Powertrain. The other attributes does not improve the accuracy score.

ROC Curve

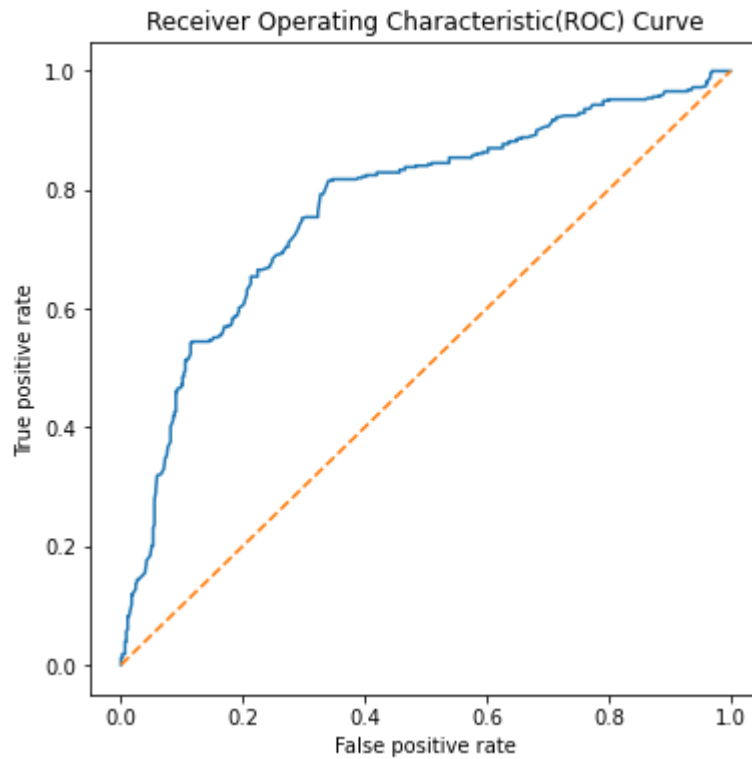
In [87]:

```
model = GaussianNB()
clf = model.fit(X_new_train, y_train)

ypredprobs = model.predict_proba(X_new_test)
probs = ypredprobs[:, 1]

fper, tper, threshold = roc_curve(ytest, probs)

plt.figure(figsize=(6, 6))
plt.plot(fper, tper)
plt.plot([0, 1], [0, 1], linestyle='--')
plt.xlabel('False positive rate')
plt.ylabel('True positive rate')
plt.title('Receiver Operating Characteristic(ROC) Curve')
plt.show()
```



This curve shows the performance of the best model that we've got with Naive Bayes method. It shows classification on the combination model between Gaussian and Categorical NB at all classification threshold. As the curve gets closer to the top-left corner, it indicates a good performance.

Model Evaluation

As a conclusion from Naive Bayes model, this method does not perform extremely good. It might be because this method really depends on the each attributes probability, their correlation between each others and the predicted value, also the prior probability (relates to the relation between predicted and another attribute). From the cross validation score, Categorical Naive Bayes got the highest score. The evaluation of the model is shown as accuracy score. Overall, the highest accuracy score is 73.15%, which is obtained by combination naive bayes method between Gaussian Naive Bayes and Categorical Naive Bayes.

7 Conclusion and Comparison

We implemented 4 prediction models for this assignment, SVM (binary classification), Neural Network (regression), Lasso (regression), and Naive Bayes using the Euro_6_latest.csv dataset to predict CO Emission.

The SVM and NN models performed well and show good results while the Lasso regression model showed very average result. SVM accuracy were 88% on the train data and 87% on the test data.

The NN model MSE were around 2% on the train data and 10%-18% on the test data on average which is very good considering the dataset features have many outliers.

The Regression Lasso did badly. The MSE were 63% for training and 67% for test data. We could not reduce the MSE even when choosing the best alpha.

The Naive Bayes model did not perform extremely good. the highest accuracy score was 73.15% which is obtained by combination naive bayes method between Gaussian Naive Bayes and Categorical Naive Bayes