

Predict the price of the Uber ride from a given pickup point to the agreed drop-off location. Perform following tasks:

1. Pre-process the data
2. Identify outliers.
3. Check the correlation.
4. Implement linear regression and random forest regression models.
5. Evaluate the models and compare their respective scores like R2, RMSE, etc. Dataset link:

<https://www.kaggle.com/datasets/yasserh/uber-fares-dataset> (<https://www.kaggle.com/datasets/yasserh/uber-fares-dataset>)

```
In [1]: #Importing the required Libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [2]: #importing the dataset
df = pd.read_csv("uber.csv")
```

1. Pre-process the dataset.

```
In [3]: df.head()
```

```
Out[3]:
```

		Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	pa:
0	24238194	2015-05-07 19:52:06.0000003		7.5	2015-05-07 19:52:06 UTC	-73.999817	40.738354	-73.999512	40.723217	
1	27835199	2009-07-17 20:04:56.0000002		7.7	2009-07-17 20:04:56 UTC	-73.994355	40.728225	-73.994710	40.750325	
2	44984355	2009-08-24 21:45:00.00000061		12.9	2009-08-24 21:45:00 UTC	-74.005043	40.740770	-73.962565	40.772647	
3	25894730	2009-06-26 08:22:21.0000001		5.3	2009-06-26 08:22:21 UTC	-73.976124	40.790844	-73.965316	40.803349	
4	17610152	2014-08-28 17:47:00.000000188		16.0	2014-08-28 17:47:00 UTC	-73.925023	40.744085	-73.973082	40.761247	

```
In [4]: df.info() #To get the required information of the dataset
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200000 entries, 0 to 199999
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
 0   Unnamed: 0        200000 non-null  int64  
 1   key              200000 non-null  object  
 2   fare_amount      200000 non-null  float64 
 3   pickup_datetime  200000 non-null  object  
 4   pickup_longitude 200000 non-null  float64 
 5   pickup_latitude   200000 non-null  float64 
 6   dropoff_longitude 199999 non-null  float64 
 7   dropoff_latitude  199999 non-null  float64 
 8   passenger_count  200000 non-null  int64  
dtypes: float64(5), int64(2), object(2)
memory usage: 13.7+ MB
```

```
In [5]: df.columns #To get number of columns in the dataset
```

```
Out[5]: Index(['Unnamed: 0', 'key', 'fare_amount', 'pickup_datetime',
               'pickup_longitude', 'pickup_latitude', 'dropoff_longitude',
               'dropoff_latitude', 'passenger_count'],
              dtype='object')
```

```
In [6]: df = df.drop(['Unnamed: 0', 'key'], axis= 1) #To drop unnamed column as it isn't required
```

```
In [7]: df.head()
```

```
Out[7]:   fare_amount      pickup_datetime  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count
0          7.5  2015-05-07 19:52:06 UTC       -73.999817        40.738354       -73.999512        40.723217             1
1          7.7  2009-07-17 20:04:56 UTC       -73.994355        40.728225       -73.994710        40.750325             1
2         12.9  2009-08-24 21:45:00 UTC       -74.005043        40.740770       -73.962565        40.772647             1
3          5.3  2009-06-26 08:22:21 UTC       -73.976124        40.790844       -73.965316        40.803349             3
4         16.0  2014-08-28 17:47:00 UTC       -73.925023        40.744085       -73.973082        40.761247             5
```

```
In [8]: df.shape #To get the total (Rows,Columns)
```

```
Out[8]: (200000, 7)
```

```
In [9]: df.dtypes #To get the type of each column
```

```
Out[9]: fare_amount      float64
pickup_datetime    object
pickup_longitude    float64
pickup_latitude     float64
dropoff_longitude   float64
dropoff_latitude    float64
passenger_count     int64
dtype: object
```

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200000 entries, 0 to 199999
Data columns (total 7 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   fare_amount        200000 non-null   float64
 1   pickup_datetime    200000 non-null   object 
 2   pickup_longitude   200000 non-null   float64
 3   pickup_latitude    200000 non-null   float64
 4   dropoff_longitude  199999 non-null   float64
 5   dropoff_latitude   199999 non-null   float64
 6   passenger_count    200000 non-null   int64  
dtypes: float64(5), int64(1), object(1)
memory usage: 10.7+ MB
```

```
In [11]: df.describe() #To get statistics of each columns
```

```
Out[11]:
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
count	200000.000000	200000.000000	200000.000000	199999.000000	199999.000000	200000.000000
mean	11.359955	-72.527638	39.935885	-72.525292	39.923890	1.684535
std	9.901776	11.437787	7.720539	13.117408	6.794829	1.385997
min	-52.000000	-1340.648410	-74.015515	-3356.666300	-881.985513	0.000000
25%	6.000000	-73.992065	40.734796	-73.991407	40.733823	1.000000
50%	8.500000	-73.981823	40.752592	-73.980093	40.753042	1.000000
75%	12.500000	-73.967154	40.767158	-73.963658	40.768001	2.000000
max	499.000000	57.418457	1644.421482	1153.572603	872.697628	208.000000

Filling Missing values

```
In [12]: df.isnull().sum()
```

```
Out[12]: fare_amount      0
pickup_datetime      0
pickup_longitude      0
pickup_latitude      0
dropoff_longitude     1
dropoff_latitude     1
passenger_count      0
dtype: int64
```

```
In [13]: df['dropoff_latitude'].fillna(value=df['dropoff_latitude'].mean(),inplace = True)
df['dropoff_longitude'].fillna(value=df['dropoff_longitude'].median(),inplace = True)
```

```
In [14]: df.isnull().sum()
```

```
Out[14]: fare_amount      0
pickup_datetime      0
pickup_longitude      0
pickup_latitude      0
dropoff_longitude     0
dropoff_latitude     0
passenger_count      0
dtype: int64
```

```
In [15]: df.dtypes
```

```
Out[15]: fare_amount      float64
pickup_datetime      object
pickup_longitude      float64
pickup_latitude      float64
dropoff_longitude     float64
dropoff_latitude     float64
passenger_count      int64
dtype: object
```

Column pickup_datetime is in wrong format (Object). Convert it to DateTime Format

```
In [16]: df.pickup_datetime = pd.to_datetime(df.pickup_datetime, errors='coerce')
```

```
In [17]: df.dtypes
```

```
Out[17]: fare_amount          float64
pickup_datetime    datetime64[ns, UTC]
pickup_longitude      float64
pickup_latitude       float64
dropoff_longitude     float64
dropoff_latitude      float64
passenger_count        int64
dtype: object
```

To segregate each time of date and time

```
In [18]: df= df.assign(hour = df.pickup_datetime.dt.hour,
                     day= df.pickup_datetime.dt.day,
                     month = df.pickup_datetime.dt.month,
                     year = df.pickup_datetime.dt.year,
                     dayofweek = df.pickup_datetime.dt.dayofweek)
```

```
In [19]: df.head()
```

```
Out[19]:   fare_amount  pickup_datetime  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count  hour  day  mon
0         7.5  2015-05-07
           19:52:06+00:00      -73.999817      40.738354     -73.999512      40.723217           1      19      7
1         7.7  2009-07-17
           20:04:56+00:00      -73.994355      40.728225     -73.994710      40.750325           1      20     17
2        12.9  2009-08-24
           21:45:00+00:00      -74.005043      40.740770     -73.962565      40.772647           1      21     24
3         5.3  2009-06-26
           08:22:21+00:00      -73.976124      40.790844     -73.965316      40.803349           3      8     26
4        16.0  2014-08-28
           17:47:00+00:00      -73.925023      40.744085     -73.973082      40.761247           5     17     28
```

```
In [20]: # drop the column 'pickup_datetime' using drop()  
# 'axis = 1' drops the specified column  
  
df = df.drop('pickup_datetime',axis=1)
```

```
In [21]: df.head()
```

```
Out[21]:   fare_amount  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count  hour  day  month  year  dayofweek  
0          7.5        -73.999817      40.738354       -73.999512      40.723217           1     19     7      5  2015  
1          7.7        -73.994355      40.728225       -73.994710      40.750325           1     20    17      7  2009  
2         12.9        -74.005043      40.740770       -73.962565      40.772647           1     21    24      8  2009  
3          5.3        -73.976124      40.790844       -73.965316      40.803349           3     8    26      6  2009  
4         16.0        -73.925023      40.744085       -73.973082      40.761247           5    17    28      8  2014
```

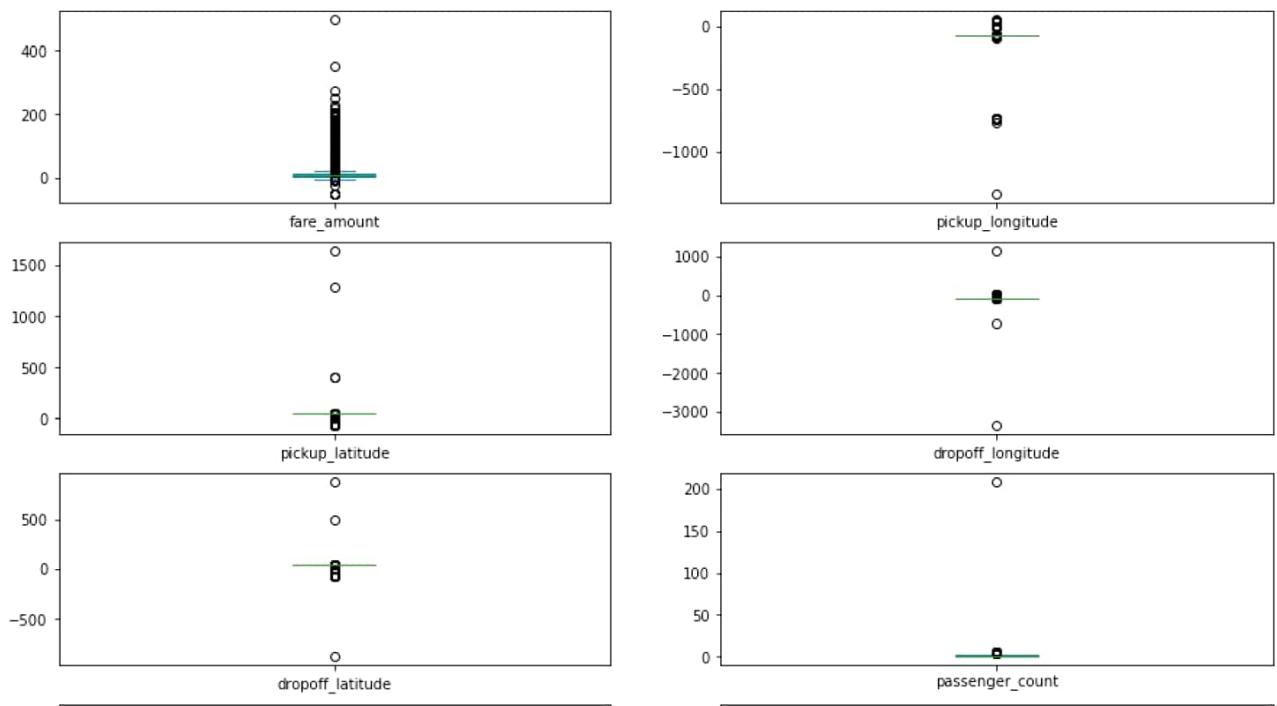
```
In [22]: df.dtypes
```

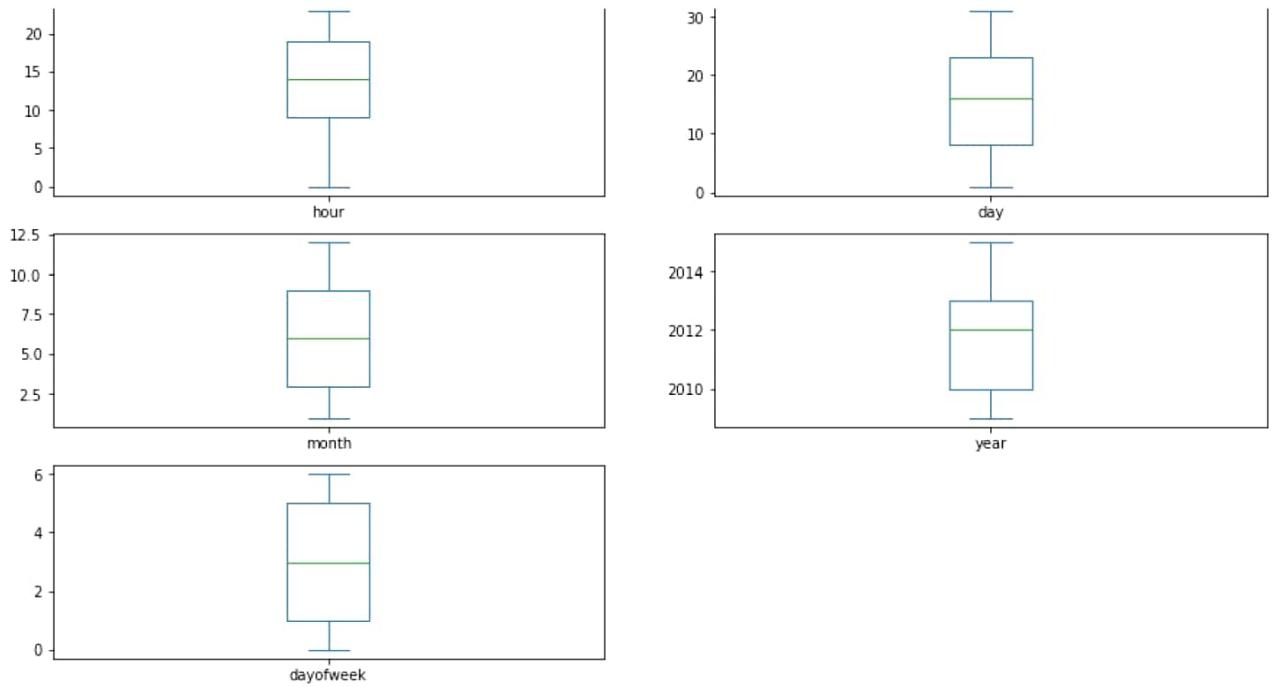
```
Out[22]: fare_amount      float64  
pickup_longitude    float64  
pickup_latitude      float64  
dropoff_longitude    float64  
dropoff_latitude      float64  
passenger_count      int64  
hour                  int64  
day                   int64  
month                 int64  
year                  int64  
dayofweek              int64  
dtype: object
```

Checking outliers and filling them

```
In [23]: df.plot(kind = "box", subplots = True, layout = (7,2), figsize=(15,20)) #Boxplot to check the outliers
```

```
Out[23]: fare_amount      AxesSubplot(0.125,0.787927;0.352273x0.0920732)
pickup_longitude    AxesSubplot(0.547727,0.787927;0.352273x0.0920732)
pickup_latitude      AxesSubplot(0.125,0.677439;0.352273x0.0920732)
dropoff_longitude   AxesSubplot(0.547727,0.677439;0.352273x0.0920732)
dropoff_latitude    AxesSubplot(0.125,0.566951;0.352273x0.0920732)
passenger_count     AxesSubplot(0.547727,0.566951;0.352273x0.0920732)
hour                AxesSubplot(0.125,0.456463;0.352273x0.0920732)
day                 AxesSubplot(0.547727,0.456463;0.352273x0.0920732)
month               AxesSubplot(0.125,0.345976;0.352273x0.0920732)
year                AxesSubplot(0.547727,0.345976;0.352273x0.0920732)
dayofweek            AxesSubplot(0.125,0.235488;0.352273x0.0920732)
dtype: object
```





In [24]: #Using the InterQuartile Range to fill the values

```

def remove_outlier(df1 , col):
    Q1 = df1[col].quantile(0.25)
    Q3 = df1[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_whisker = Q1-1.5*IQR
    upper_whisker = Q3+1.5*IQR
    df[col] = np.clip(df1[col] , lower_whisker , upper_whisker)
    return df1

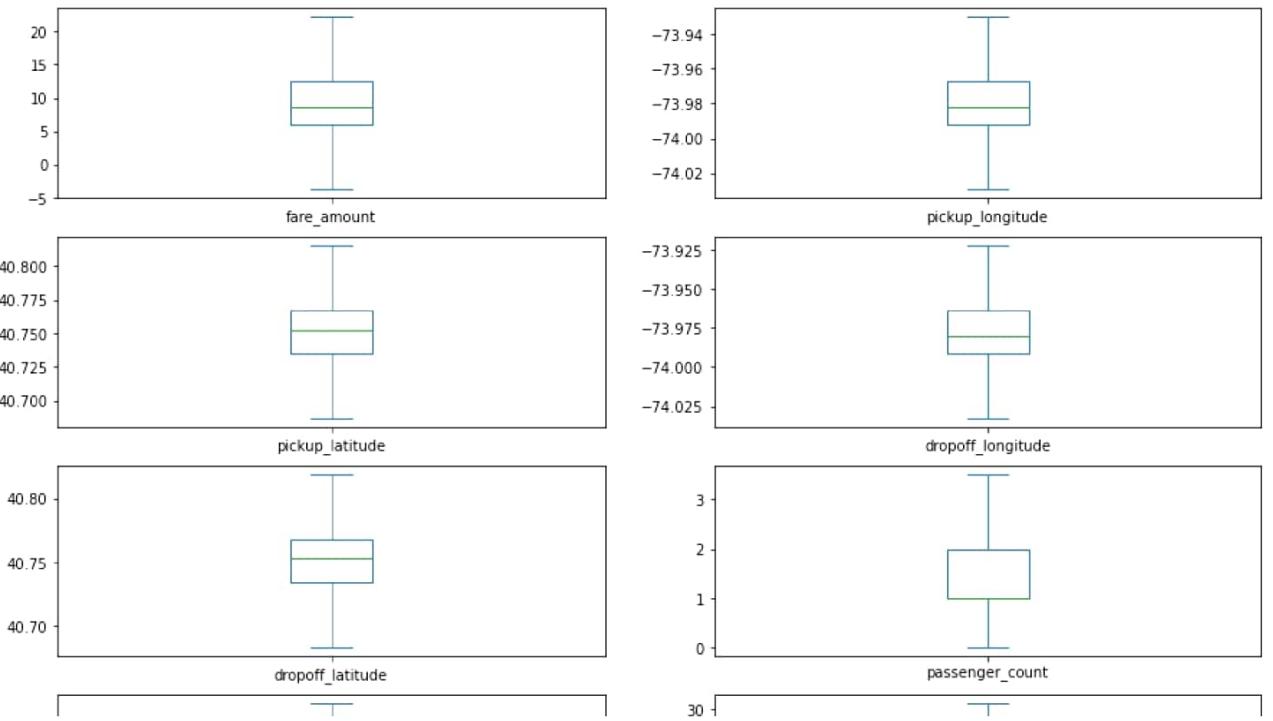
def treat_outliers_all(df1 , col_list):
    for c in col_list:
        df1 = remove_outlier(df1 , c)
    return df1

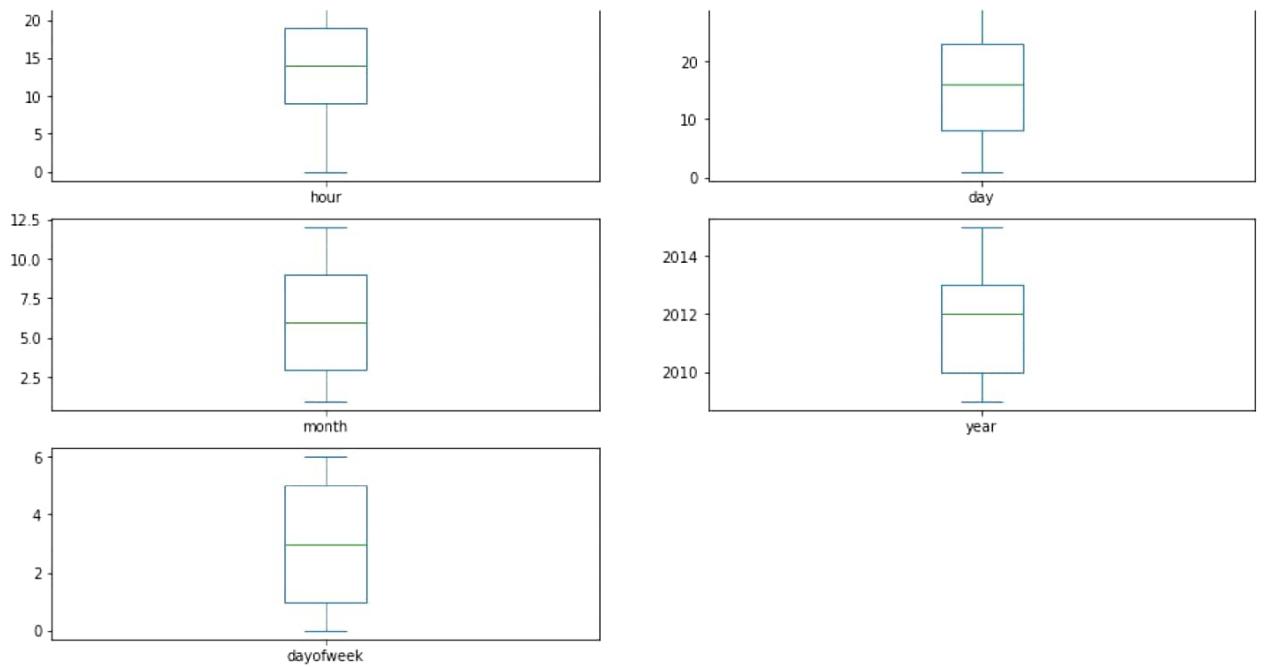
```

```
In [25]: df = treat_outliers_all(df , df.iloc[:, 0::])
```

```
In [26]: df.plot(kind = "box", subplots = True, layout = (7,2), figsize=(15,20)) #Boxplot shows that dataset is free from outliers
```

```
Out[26]: fare_amount      AxesSubplot(0.125,0.787927;0.352273x0.0920732)
pickup_longitude    AxesSubplot(0.547727,0.787927;0.352273x0.0920732)
pickup_latitude       AxesSubplot(0.125,0.677439;0.352273x0.0920732)
dropoff_longitude    AxesSubplot(0.547727,0.677439;0.352273x0.0920732)
dropoff_latitude      AxesSubplot(0.125,0.566951;0.352273x0.0920732)
passenger_count       AxesSubplot(0.547727,0.566951;0.352273x0.0920732)
hour                  AxesSubplot(0.125,0.456463;0.352273x0.0920732)
day                   AxesSubplot(0.547727,0.456463;0.352273x0.0920732)
month                 AxesSubplot(0.125,0.345976;0.352273x0.0920732)
year                  AxesSubplot(0.547727,0.345976;0.352273x0.0920732)
dayofweek              AxesSubplot(0.125,0.235488;0.352273x0.0920732)
dtype: object
```





```
In [28]: pip install haversine
```

```
Collecting haversine
  Note: you may need to restart the kernel to use updated packages.
    Downloading haversine-2.6.0-py2.py3-none-any.whl (6.8 kB)
Installing collected packages: haversine
Successfully installed haversine-2.6.0
```

```
In [30]: #pip install haversine
import haversine as hs  #Calculate the distance using Haversine to calculate the distance between two points. Can
travel_dist = []
for pos in range(len(df['pickup_longitude'])):
    long1,lat1,long2,lat2 = [df['pickup_longitude'][pos],df['pickup_latitude'][pos],df['dropoff_longitude'][pos],df['dropoff_latitude'][pos]]
    loc1=(lat1,long1)
    loc2=(lat2,long2)
    c = hs.haversine(loc1,loc2)
    travel_dist.append(c)

print(travel_dist)
df['dist_travel_km'] = travel_dist
df.head()
```

IOPub data rate exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub_data_rate_limit`.

Current values:
NotebookApp.iopub_data_rate_limit=1000000.0 (bytes/sec)
NotebookApp.rate_limit_window=3.0 (secs)

```
Out[30]:   fare_amount  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count  hour  day  month  year  dayofweek
0          7.5        -73.999817      40.738354       -73.999512      40.723217           1.0     19     7      5  2015
1          7.7        -73.994355      40.728225       -73.994710      40.750325           1.0     20    17      7  2009
2         12.9        -74.005043      40.740770       -73.962565      40.772647           1.0     21    24      8  2009
3          5.3        -73.976124      40.790844       -73.965316      40.803349           3.0     8    26      6  2009
4         16.0        -73.929786      40.744085       -73.973082      40.761247           3.5    17    28      8  2014
```

```
In [31]: #Uber doesn't travel over 130 kms so minimize the distance
df= df.loc[(df.dist_travel_km >= 1) | (df.dist_travel_km <= 130)]
print("Remaining observations in the dataset:", df.shape)
```

Remaining observations in the dataset: (200000, 12)

```
In [32]: #Finding incorrect Latitude (Less than or greater than 90) and Longitude (greater than or Less than 180)
incorrect_coordinates = df.loc[(df.pickup_latitude > 90) |(df.pickup_latitude < -90) |
                               (df.dropoff_latitude > 90) |(df.dropoff_latitude < -90) |
                               (df.pickup_longitude > 180) |(df.pickup_longitude < -180) |
                               (df.dropoff_longitude > 90) |(df.dropoff_longitude < -90)
                               ]
```

```
In [33]: df.drop(incorrect_coordinates, inplace = True, errors = 'ignore')
```

```
In [34]: df.head()
```

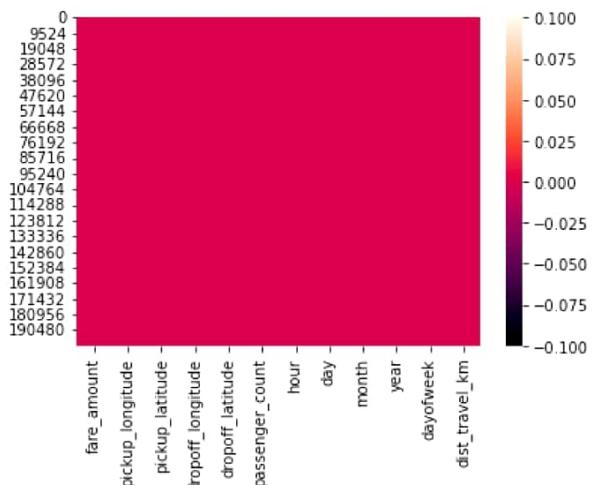
```
Out[34]:   fare_amount  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count  hour  day  month  year  dayofweek
0          7.5      -73.999817     40.738354      -73.999512      40.723217           1.0    19    7    5  2015
1          7.7      -73.994355     40.728225      -73.994710      40.750325           1.0    20   17    7  2009
2         12.9      -74.005043     40.740770      -73.962565      40.772647           1.0    21   24    8  2009
3          5.3      -73.976124     40.790844      -73.965316      40.803349           3.0     8   26    6  2009
4         16.0      -73.929786     40.744085      -73.973082      40.761247           3.5    17   28    8  2014
```

```
In [35]: df.isnull().sum()
```

```
Out[35]: fare_amount      0
pickup_longitude     0
pickup_latitude       0
dropoff_longitude     0
dropoff_latitude       0
passenger_count       0
hour                   0
day                     0
month                  0
year                   0
dayofweek                0
dist_travel_km         0
dtype: int64
```

```
In [36]: sns.heatmap(df.isnull()) #Free for null values
```

```
Out[36]: <AxesSubplot:>
```



```
In [37]: corr = df.corr() #Function to find the correlation
```

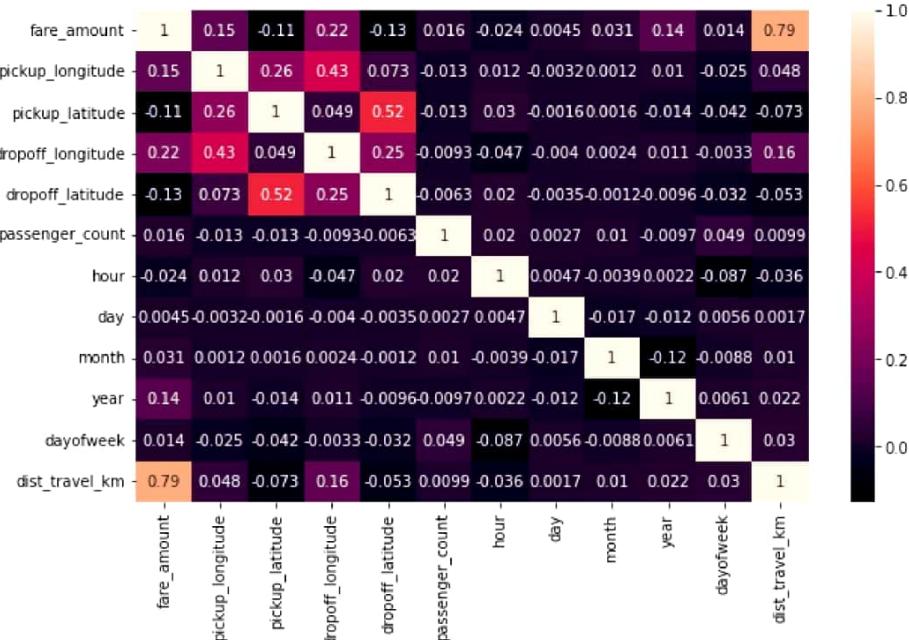
```
In [38]: corr
```

```
Out[38]:
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day
fare_amount	1.000000	0.154069	-0.110842	0.218675	-0.125898	0.015778	-0.023623	0.0045
pickup_longitude	0.154069	1.000000	0.259497	0.425619	0.073290	-0.013213	0.011579	-0.0032
pickup_latitude	-0.110842	0.259497	1.000000	0.048889	0.515714	-0.012889	0.029681	-0.0015
dropoff_longitude	0.218675	0.425619	0.048889	1.000000	0.245667	-0.009303	-0.046558	-0.0040
dropoff_latitude	-0.125898	0.073290	0.515714	0.245667	1.000000	-0.006308	0.019783	-0.0034
passenger_count	0.015778	-0.013213	-0.012889	-0.009303	-0.006308	1.000000	0.020274	0.0027
hour	-0.023623	0.011579	0.029681	-0.046558	0.019783	0.020274	1.000000	0.0046
day	0.004534	-0.003204	-0.001553	-0.004007	-0.003479	0.002712	0.004677	1.0000
month	0.030817	0.001169	0.001562	0.002391	-0.001193	0.010351	-0.003926	-0.0173
year	0.141277	0.010198	-0.014243	0.011346	-0.009603	-0.009749	0.002156	-0.0121
dayofweek	0.013652	-0.024652	-0.042310	-0.003336	-0.031919	0.048550	-0.086947	0.0056
dist_travel_km	0.786385	0.048446	-0.073362	0.155191	-0.052701	0.009884	-0.035708	0.0017

```
In [39]: fig,ax = plt.subplots(figsize = (10,6))
sns.heatmap(df.corr(),annot = True) #Correlation Heatmap (Light values means highly correlated)
```

```
Out[39]: <AxesSubplot:>
```



Dividing the dataset into feature and target values

```
In [40]: x = df[['pickup_longitude','pickup_latitude','dropoff_longitude','dropoff_latitude','passenger_count','hour','da  
↳ ↴  
In [41]: y = df['fare_amount']
```

Dividing the dataset into training and testing dataset

```
In [42]: from sklearn.model_selection import train_test_split  
X_train,X_test,y_train,y_test = train_test_split(x,y,test_size = 0.33)
```

Linear Regression

```
In [43]: from sklearn.linear_model import LinearRegression  
regression = LinearRegression()  
  
In [44]: regression.fit(X_train,y_train)  
Out[44]: LinearRegression()  
  
In [45]: regression.intercept_ #To find the Linear intercept  
Out[45]: 3727.0995043526045  
  
In [46]: regression.coef_ #To find the Linear coefficient  
Out[46]: array([ 2.56317588e+01, -7.16931941e+00,  2.05592274e+01, -1.83436182e+01,  
    7.56526199e-02,  5.17343471e-03,  2.74067422e-03,  5.95389786e-02,  
    3.64889506e-01, -3.15727784e-02,  1.84522518e+00])
```

```
In [47]: prediction = regression.predict(X_test) #To predict the target values
```

```
In [48]: print(prediction)
```

```
[12.6209887  4.39051806  6.46240798 ... 33.45494174  7.00560759  
 15.70280297]
```

```
In [49]: y_test
```

```
Out[49]: 154908    9.50  
      5277    4.50  
     170208    6.10  
     8592    12.00  
    83444    5.00  
      ...  
    99421   22.25  
   83027    6.00  
   26526   22.25  
  13909    4.90  
  42270   20.50  
Name: fare_amount, Length: 66000, dtype: float64
```

Metrics Evaluation using R2, Mean Squared Error, Root Mean Squared Error

```
In [50]: from sklearn.metrics import r2_score
```

```
In [51]: r2_score(y_test,prediction)
```

```
Out[51]: 0.661012080097631
```

```
In [52]: from sklearn.metrics import mean_squared_error
```

```
In [53]: MSE = mean_squared_error(y_test,prediction)
```

```
In [54]: MSE  
Out[54]: 10.093810949414875  
  
In [55]: RMSE = np.sqrt(MSE)  
  
In [56]: RMSE  
Out[56]: 3.1770758488608477
```

Random Forest Regression

```
In [60]: from sklearn.ensemble import RandomForestRegressor  
  
In [61]: rf = RandomForestRegressor(n_estimators=100) #Here n_estimators means number of trees you want to build before m  
        ↪  
In [62]: rf.fit(X_train,y_train)  
Out[62]: RandomForestRegressor()  
  
In [63]: y_pred = rf.predict(X_test)  
  
In [64]: y_pred  
Out[64]: array([12.5475, 4.391, 7.028, ..., 21.155, 5.919, 16.379])
```

Metrics evaluatin for Random Forest

```
In [65]: R2_Random = r2_score(y_test,y_pred)
```

```
In [66]: R2_Random
```

```
Out[66]: 0.7937449132954233
```

```
In [67]: MSE_Random = mean_squared_error(y_test,y_pred)
```

```
In [68]: MSE_Random
```

```
Out[68]: 6.141516350053928
```

```
In [69]: RMSE_Random = np.sqrt(MSE_Random)
```

```
In [70]: RMSE_Random
```

```
Out[70]: 2.4782082943235277
```

Assignment 2

2. Classify the email using the binary classification method. Email Spam detection has two states: a) Normal State – Not Spam, b) Abnormal State – Spam. Use K-Nearest Neighbors and Support Vector Machine for classification. Analyze their performance. Dataset link: The emails.csv dataset on the Kaggle <https://www.kaggle.com/datasets/balaka18/email-spam-classification-dataset-csv> (<https://www.kaggle.com/datasets/balaka18/email-spam-classification-dataset-csv>)

```
In [3]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn import metrics
```

```
In [4]: df=pd.read_csv('emails.csv')
```

```
In [5]: df.head()
```

Out[5]:

Email No.	the	to	ect	and	for	of	a	you	hou	...	connevey	jay	valued	lay	infrastructure	
0 Email 1	0	0	1	0	0	0	2	0	0	...	0	0	0	0	0	C
1 Email 2	8	13	24	6	6	2	102	1	27	...	0	0	0	0	0	C
2 Email 3	0	0	1	0	0	0	8	0	0	...	0	0	0	0	0	C
3 Email 4	0	5	22	0	5	1	51	2	10	...	0	0	0	0	0	C
4 Email 5	7	6	17	1	5	2	57	0	9	...	0	0	0	0	0	C

5 rows × 3002 columns



```
In [7]: df.columns
```

```
Out[7]: Index(['Email No.', 'the', 'to', 'ect', 'and', 'for', 'of', 'a', 'you', 'hou',
               ...
               'connevey', 'jay', 'valued', 'lay', 'infrastructure', 'military',
               'allowing', 'ff', 'dry', 'Prediction'],
               dtype='object', length=3002)
```

```
In [8]: df.isnull().sum()
```

```
Out[8]: Email No.      0
the          0
to           0
ect          0
and          0
...
military     0
allowing     0
ff            0
dry           0
Prediction    0
Length: 3002, dtype: int64
```

```
In [9]: df.dropna(inplace = True)
```

```
In [10]: df.drop(['Email No.'],axis=1,inplace=True)
X = df.drop(['Prediction'],axis = 1)
y = df['Prediction']
```

```
In [11]: from sklearn.preprocessing import scale
X = scale(X)
# split into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 42)
```

##KNN classifier

```
In [12]: from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=7)

knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
```

```
In [13]: print("Prediction",y_pred)
```

```
Prediction [0 0 1 ... 1 1 1]
```

```
In [14]: print("KNN accuracy = ",metrics.accuracy_score(y_test,y_pred))
```

```
KNN accuracy =  0.8009020618556701
```

```
In [15]: print("Confusion matrix",metrics.confusion_matrix(y_test,y_pred))
```

```
Confusion matrix [[804 293]
 [ 16 439]]
```

SVM classifier

```
In [18]: # cost C = 1
model = SVC(C = 1)

# fit
model.fit(X_train, y_train)

# predict
y_pred = model.predict(X_test)
```

```
In [19]: metrics.confusion_matrix(y_true=y_test, y_pred=y_pred)
```

```
Out[19]: array([[1091,    6],
                 [  90, 365]], dtype=int64)
```

```
In [20]: print("SVM accuracy = ",metrics.accuracy_score(y_test,y_pred))
```

```
SVM accuracy =  0.9381443298969072
```

Assignment 5

KNN algorithm on diabetes dataset

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn import metrics
```

```
In [2]: df=pd.read_csv('diabetes.csv')
```

```
In [3]: df.columns
```

```
Out[3]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'Pedigree', 'Age', 'Outcome'],
       dtype='object')
```

Check for null values. If present remove null values from the dataset

```
In [4]: df.isnull().sum()
```

```
Out[4]: Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
Pedigree         0
Age              0
Outcome          0
dtype: int64
```

Outcome is the label/target, other columns are features

```
In [5]: X = df.drop('Outcome',axis = 1)
y = df['Outcome']
```

```
In [6]: from sklearn.preprocessing import scale
X = scale(X)
# split into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random
```

```
In [7]: from sklearn.neighbors import KNeighborsClassifier  
knn = KNeighborsClassifier(n_neighbors=7)  
  
knn.fit(X_train, y_train)  
y_pred = knn.predict(X_test)
```

```
In [8]: print("Confusion matrix: ")  
cs = metrics.confusion_matrix(y_test,y_pred)  
print(cs)
```

```
Confusion matrix:  
[[123  28]  
 [ 37  43]]
```

```
In [9]: print("Accuracy ",metrics.accuracy_score(y_test,y_pred))
```

```
Accuracy 0.7186147186147186
```

Classification error rate: proportion of instances misclassified over the whole set of instances. Error rate is calculated as the total number of two incorrect predictions (FN + FP) divided by the total number of a dataset (examples in the dataset).

Also error_rate = 1- accuracy

```
In [10]: total_misclassified = cs[0,1] + cs[1,0]  
print(total_misclassified)  
total_examples = cs[0,0]+cs[0,1]+cs[1,0]+cs[1,1]  
print(total_examples)  
print("Error rate",total_misclassified/total_examples)  
print("Error rate ",1-metrics.accuracy_score(y_test,y_pred))
```

```
65  
231  
Error rate 0.2813852813852814  
Error rate 0.2813852813852814
```

```
In [11]: print("Precision score",metrics.precision_score(y_test,y_pred))
```

```
Precision score 0.6056338028169014
```

```
In [12]: print("Recall score ",metrics.recall_score(y_test,y_pred))
```

```
Recall score 0.5375
```

```
In [13]: print("Classification report ",metrics.classification_report(y_test,y_pred))
```

Classification report		precision	recall	f1-score	support
0	0.77	0.81	0.79	151	
1	0.61	0.54	0.57	80	
accuracy			0.72	231	
macro avg	0.69	0.68	0.68	231	
weighted avg	0.71	0.72	0.71	231	

Given a bank customer, build a neural network-based classifier that can determine whether they will leave or not in the next 6 months.

Dataset Description: The case study is from an open-source dataset from Kaggle. The dataset contains 10,000 sample points with 14 distinct features such as CustomerId, CreditScore, Geography, Gender, Age, Tenure, Balance, etc. Link to the Kaggle project:

<https://www.kaggle.com/barelydedicated/bank-customer-churn-modeling> Perform following steps:

1. Read the dataset.
2. Distinguish the feature and target set and divide the data set into training and test sets.
3. Normalize the train and test data.
4. Initialize and build the model. Identify the points of improvement and implement the same.
5. Print the accuracy score and confusion matrix.

In [1]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt #Importing the Libraries
```

In [2]:

```
df = pd.read_csv("Churn_Modelling.csv")
```

Preprocessing.

In [3]:

```
df.head()
```

Out[3]:

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	I
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	
2	3	15619304	Onio	502	France	Female	42	8	159660.80	
3	4	15701354	Boni	699	France	Female	39	1	0.00	
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	



In [4]:

```
df.shape
```

Out[4]:

```
(10000, 14)
```

In [5]:

```
df.describe()
```

Out[5]:

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumO
count	10000.00000	1.000000e+04	10000.00000	10000.00000	10000.00000	10000.00000	100
mean	5000.50000	1.569094e+07	650.528800	38.921800	5.012800	76485.889288	
std	2886.89568	7.193619e+04	96.653299	10.487806	2.892174	62397.405202	
min	1.00000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	
25%	2500.75000	1.562853e+07	584.000000	32.000000	3.000000	0.000000	
50%	5000.50000	1.569074e+07	652.000000	37.000000	5.000000	97198.540000	
75%	7500.25000	1.575323e+07	718.000000	44.000000	7.000000	127644.240000	
max	10000.00000	1.581569e+07	850.000000	92.000000	10.000000	250898.090000	

◀ ▶

In [6]:

df.isnull()

Out[6]:

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...
9995	False	False	False	False	False	False	False	False	False
9996	False	False	False	False	False	False	False	False	False
9997	False	False	False	False	False	False	False	False	False
9998	False	False	False	False	False	False	False	False	False
9999	False	False	False	False	False	False	False	False	False

10000 rows × 14 columns

◀ ▶

In [7]:

df.isnull().sum()

Out[7]:

```
RowNumber      0
CustomerId     0
Surname        0
CreditScore    0
Geography      0
Gender         0
Age            0
Tenure         0
Balance        0
NumOfProducts  0
HasCrCard      0
IsActiveMember 0
EstimatedSalary 0
Exited         0
dtype: int64
```

```
In [8]:
```

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   RowNumber        10000 non-null   int64  
 1   CustomerId       10000 non-null   int64  
 2   Surname          10000 non-null   object  
 3   CreditScore      10000 non-null   int64  
 4   Geography         Geography        object  
 5   Gender           10000 non-null   object  
 6   Age              10000 non-null   int64  
 7   Tenure           10000 non-null   int64  
 8   Balance          10000 non-null   float64 
 9   NumOfProducts    10000 non-null   int64  
 10  HasCrCard        10000 non-null   int64  
 11  IsActiveMember   10000 non-null   int64  
 12  EstimatedSalary  10000 non-null   float64 
 13  Exited          10000 non-null   int64  
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

```
In [9]:
```

```
df.dtypes
```

```
Out[9]:
```

```
RowNumber        int64
CustomerId       int64
Surname          object
CreditScore      int64
Geography        object
Gender           object
Age              int64
Tenure           int64
Balance          float64
NumOfProducts    int64
HasCrCard        int64
IsActiveMember   int64
EstimatedSalary  float64
Exited          int64
dtype: object
```

```
In [10]:
```

```
df.columns
```

```
Out[10]:
```

```
Index(['RowNumber', 'CustomerId', 'Surname', 'CreditScore', 'Geography',
       'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard',
       'IsActiveMember', 'EstimatedSalary', 'Exited'],
      dtype='object')
```

```
In [11]:
```

```
df = df.drop(['RowNumber', 'Surname', 'CustomerId'], axis=1) #Dropping the unnecessary columns
```

```
In [12]:
```

```
df.head()
```

```
Out[12]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember
0	619	France	Female	42	2	0.00		1	1
1	608	Spain	Female	41	1	83807.86		1	0
2	502	France	Female	42	8	159660.80		3	1
3	699	France	Female	39	1	0.00		2	0

CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveN
4	850	Spain	Female	43	2	125510.82	1	1

Visualization

In [13]:

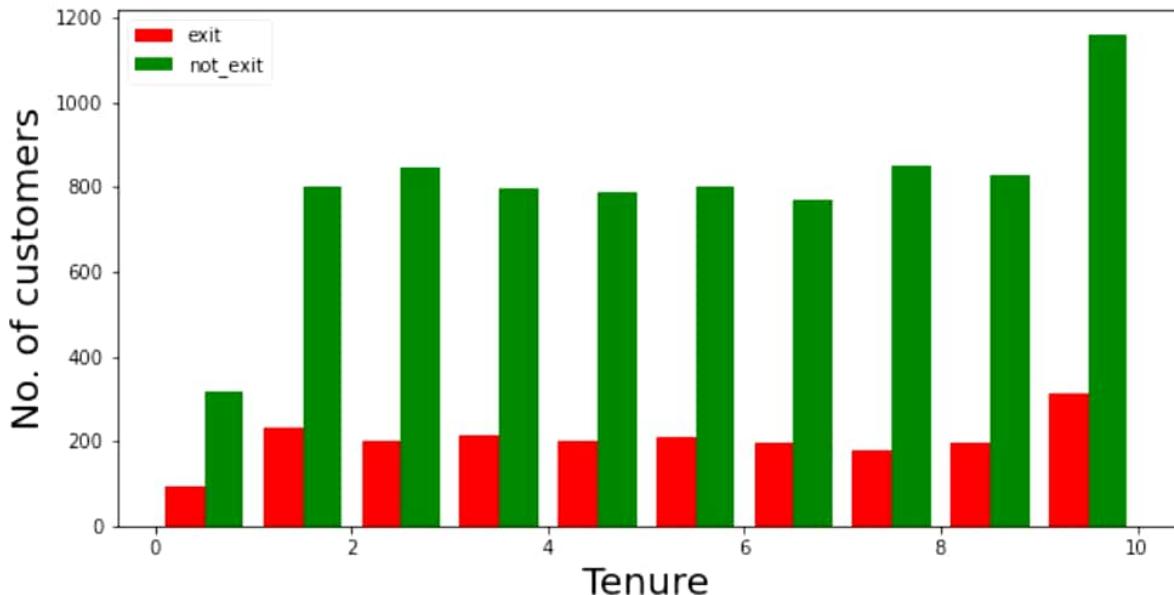
```
def visualization(x, y, xlabel):
    plt.figure(figsize=(10,5))
    plt.hist([x, y], color=['red', 'green'], label = ['exit', 'not_exit'])
    plt.xlabel(xlabel, fontsize=20)
    plt.ylabel("No. of customers", fontsize=20)
    plt.legend()
```

In [14]:

```
df_churn_exited = df[df['Exited']==1]['Tenure']
df_churn_not_exited = df[df['Exited']==0]['Tenure']
```

In [15]:

```
visualization(df_churn_exited, df_churn_not_exited, "Tenure")
```

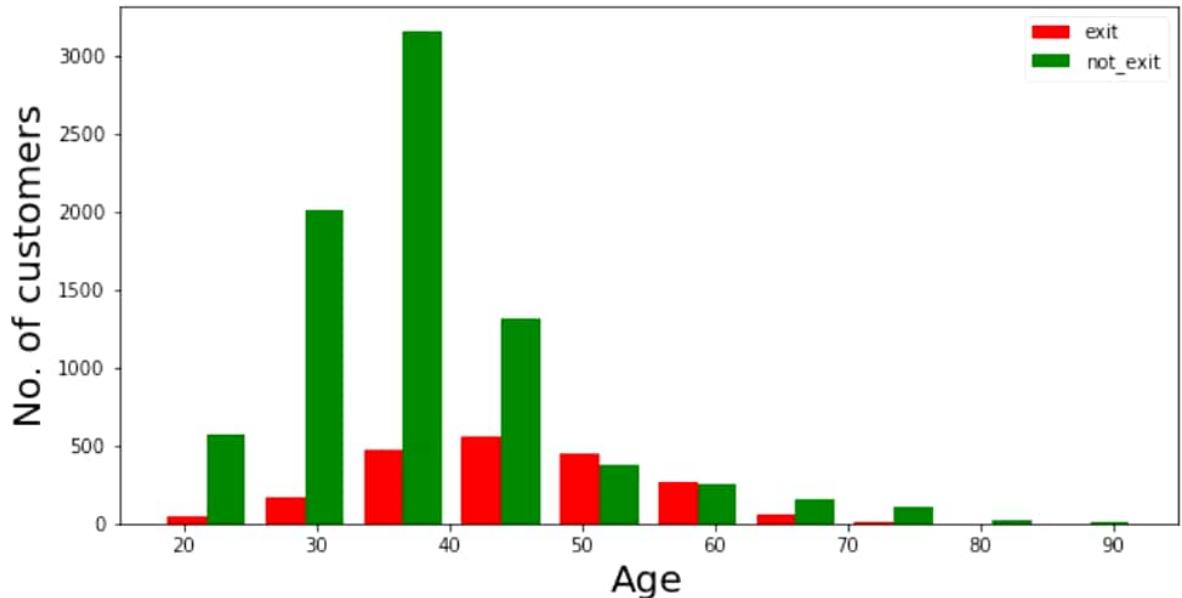


In [16]:

```
df_churn_exited2 = df[df['Exited']==1]['Age']
df_churn_not_exited2 = df[df['Exited']==0]['Age']
```

In [17]:

```
visualization(df_churn_exited2, df_churn_not_exited2, "Age")
```



Converting the Categorical Variables

```
In [18]: X = df[['CreditScore', 'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'states = pd.get_dummies(df['Geography'], drop_first = True)
gender = pd.get_dummies(df['Gender'], drop_first = True)]
```

```
In [19]: df = pd.concat([df, gender, states], axis = 1)
```

Splitting the training and testing Dataset

```
In [20]: df.head()
```

```
Out[20]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveN
0	619	France	Female	42	2	0.00		1	1
1	608	Spain	Female	41	1	83807.86		1	0
2	502	France	Female	42	8	159660.80		3	1
3	699	France	Female	39	1	0.00		2	0
4	850	Spain	Female	43	2	125510.82		1	1

```
In [21]: X = df[['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveN]]
```

```
In [22]: y = df['Exited']
```

```
In [23]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30)
```

Normalizing the values with mean as 0 and Standard Deviation as 1

```
In [24]: from sklearn.preprocessing import StandardScaler  
sc = StandardScaler()
```

```
In [25]: X_train = sc.fit_transform(X_train)  
X_test = sc.transform(X_test)
```

```
In [26]: X_train
```

```
Out[26]: array([[ 0.43051043, -0.84641337,  0.34043863, ..., -1.08909764,  
   -0.57932976, -0.57295135],  
   [-1.51285658, -0.27575181, -1.0424544 , ..., -1.08909764,  
   -0.57932976, -0.57295135],  
   [-0.71690307, -0.65619285, -1.38817765, ..., -1.08909764,  
   -0.57932976, -0.57295135],  
   ...,  
   [ 0.22376926,  0.19979948,  1.3776084 , ..., -1.08909764,  
   1.72613262, -0.57295135],  
   [ 1.92938392, -0.37086207, -1.73390091, ...,  0.91819132,  
   -0.57932976,  1.7453489 ],  
   [-0.98566659,  1.15090206,  1.03188514, ..., -1.08909764,  
   -0.57932976, -0.57295135]])
```

```
In [27]: X_test
```

```
Out[27]: array([[ 0.27545455, -0.27575181, -1.38817765, ...,  0.91819132,  
   1.72613262, -0.57295135],  
   [ 0.92668924,  0.00957896, -0.00528463, ...,  0.91819132,  
   -0.57932976, -0.57295135],  
   [ 1.44354217,  1.24601232,  0.68616189, ..., -1.08909764,  
   -0.57932976, -0.57295135],  
   ...,  
   [ 2.05342863,  0.86557129,  0.34043863, ..., -1.08909764,  
   1.72613262, -0.57295135],  
   [-1.15105953,  0.00957896, -1.0424544 , ...,  0.91819132,  
   -0.57932976,  1.7453489 ],  
   [ 2.05342863, -0.65619285, -0.35100788, ...,  0.91819132,  
   -0.57932976, -0.57295135]])
```

Building the Classifier Model using Keras

```
In [28]: import keras  
#Keras is the wrapper on the top of tensorflow  
#Can use Tensorflow as well but won't be able to understand the errors initially.
```

```
In [29]: from keras.models import Sequential #To create sequential neural network  
from keras.layers import Dense #To create hidden layers
```

```
In [30]: classifier = Sequential()
```

```
In [31]: #To add the Layers
```

```

#Dense helps to construct the neurons
#Input Dimension means we have 11 features
# Units is to create the hidden layers
#Uniform helps to distribute the weight uniformly
classifier.add(Dense(activation = "relu",input_dim = 11,units = 6,kernel_initializer = "uniform"))

In [32]: classifier.add(Dense(activation = "relu",units = 6,kernel_initializer = "uniform"))

In [33]: classifier.add(Dense(activation = "sigmoid",units = 1,kernel_initializer = "uniform"))

In [34]: classifier.compile(optimizer="adam",loss = 'binary_crossentropy',metrics = ['accuracy'])

In [35]: classifier.summary() #3 Layers created. 6 neurons in 1st,6neurons in 2nd Layer and 1 neuron in output layer

Model: "sequential"



| Layer (type)    | Output Shape | Param # |
|-----------------|--------------|---------|
| dense (Dense)   | (None, 6)    | 72      |
| dense_1 (Dense) | (None, 6)    | 42      |
| dense_2 (Dense) | (None, 1)    | 7       |


Total params: 121
Trainable params: 121
Non-trainable params: 0

In [36]: classifier.fit(X_train,y_train,batch_size=10,epochs=50) #Fitting the ANN to training set

Epoch 1/50
700/700 [=====] - 1s 1ms/step - loss: 0.4851 - accuracy: 0.7953
Epoch 2/50
700/700 [=====] - 1s 1ms/step - loss: 0.4116 - accuracy: 0.8241
Epoch 3/50
700/700 [=====] - 1s 1ms/step - loss: 0.3948 - accuracy: 0.8309
Epoch 4/50
700/700 [=====] - 1s 978us/step - loss: 0.3853 - accuracy: 0.8331
Epoch 5/50
700/700 [=====] - 1s 922us/step - loss: 0.3786 - accuracy: 0.8364
Epoch 6/50
700/700 [=====] - 1s 932us/step - loss: 0.3741 - accuracy: 0.8439
Epoch 7/50
700/700 [=====] - 1s 962us/step - loss: 0.3698 - accuracy: 0.8483
Epoch 8/50
700/700 [=====] - 1s 854us/step - loss: 0.3672 - accuracy: 0.8491
Epoch 9/50
700/700 [=====] - 1s 826us/step - loss: 0.3646 - accuracy: 0.8499
Epoch 10/50
700/700 [=====] - 1s 1ms/step - loss: 0.3620 - accuracy: 0.8536

```

```
Epoch 11/50
700/700 [=====] - 1s 1ms/step - loss: 0.3595 - accuracy: 0.
8547
Epoch 12/50
700/700 [=====] - 1s 1ms/step - loss: 0.3577 - accuracy: 0.
8546
Epoch 13/50
700/700 [=====] - 1s 1ms/step - loss: 0.3569 - accuracy: 0.
8559
Epoch 14/50
700/700 [=====] - 1s 1ms/step - loss: 0.3550 - accuracy: 0.
8557
Epoch 15/50
700/700 [=====] - 1s 1ms/step - loss: 0.3538 - accuracy: 0.
8583
Epoch 16/50
700/700 [=====] - 1s 1ms/step - loss: 0.3526 - accuracy: 0.
8569
Epoch 17/50
700/700 [=====] - 1s 1ms/step - loss: 0.3517 - accuracy: 0.
8604
Epoch 18/50
700/700 [=====] - 1s 1ms/step - loss: 0.3520 - accuracy: 0.
8570
Epoch 19/50
700/700 [=====] - 1s 1ms/step - loss: 0.3499 - accuracy: 0.
8573
Epoch 20/50
700/700 [=====] - 1s 1ms/step - loss: 0.3490 - accuracy: 0.
8587
Epoch 21/50
700/700 [=====] - 1s 954us/step - loss: 0.3502 - accuracy:
0.8586
Epoch 22/50
700/700 [=====] - 1s 1ms/step - loss: 0.3489 - accuracy: 0.
8586
Epoch 23/50
700/700 [=====] - 1s 932us/step - loss: 0.3496 - accuracy:
0.8601
Epoch 24/50
700/700 [=====] - 1s 1ms/step - loss: 0.3487 - accuracy: 0.
8589
Epoch 25/50
700/700 [=====] - 1s 1ms/step - loss: 0.3483 - accuracy: 0.
8569
Epoch 26/50
700/700 [=====] - 1s 1ms/step - loss: 0.3475 - accuracy: 0.
8594
Epoch 27/50
700/700 [=====] - 1s 1ms/step - loss: 0.3472 - accuracy: 0.
8563
Epoch 28/50
700/700 [=====] - 1s 1ms/step - loss: 0.3450 - accuracy: 0.
8589
Epoch 29/50
700/700 [=====] - 1s 1ms/step - loss: 0.3444 - accuracy: 0.
8600
Epoch 30/50
700/700 [=====] - 1s 862us/step - loss: 0.3441 - accuracy:
0.8590
Epoch 31/50
700/700 [=====] - 1s 1ms/step - loss: 0.3437 - accuracy: 0.
8590
Epoch 32/50
700/700 [=====] - 1s 1ms/step - loss: 0.3429 - accuracy: 0.
8574
Epoch 33/50
700/700 [=====] - 1s 1ms/step - loss: 0.3431 - accuracy: 0.
8601
```

```
Epoch 34/50
700/700 [=====] - 1s 1ms/step - loss: 0.3428 - accuracy: 0.8581
Epoch 35/50
700/700 [=====] - 1s 1ms/step - loss: 0.3408 - accuracy: 0.8596
Epoch 36/50
700/700 [=====] - 1s 1ms/step - loss: 0.3420 - accuracy: 0.8597
Epoch 37/50
700/700 [=====] - 1s 969us/step - loss: 0.3418 - accuracy: 0.8596
Epoch 38/50
700/700 [=====] - 1s 907us/step - loss: 0.3408 - accuracy: 0.8597
Epoch 39/50
700/700 [=====] - 1s 982us/step - loss: 0.3405 - accuracy: 0.8584
Epoch 40/50
700/700 [=====] - 1s 1ms/step - loss: 0.3407 - accuracy: 0.8597
Epoch 41/50
700/700 [=====] - 1s 994us/step - loss: 0.3398 - accuracy: 0.8610
Epoch 42/50
700/700 [=====] - 1s 1ms/step - loss: 0.3392 - accuracy: 0.8591
Epoch 43/50
700/700 [=====] - 1s 1ms/step - loss: 0.3400 - accuracy: 0.8591
Epoch 44/50
700/700 [=====] - 1s 1ms/step - loss: 0.3393 - accuracy: 0.8607
Epoch 45/50
700/700 [=====] - 1s 988us/step - loss: 0.3392 - accuracy: 0.8581
Epoch 46/50
700/700 [=====] - 1s 1ms/step - loss: 0.3385 - accuracy: 0.8616
Epoch 47/50
700/700 [=====] - 1s 1ms/step - loss: 0.3390 - accuracy: 0.8617
Epoch 48/50
700/700 [=====] - 1s 1ms/step - loss: 0.3390 - accuracy: 0.8601
Epoch 49/50
700/700 [=====] - 1s 1ms/step - loss: 0.3394 - accuracy: 0.8614
Epoch 50/50
700/700 [=====] - 1s 992us/step - loss: 0.3380 - accuracy: 0.8604
Out[36]: <tensorflow.python.keras.callbacks.History at 0x1e4c3ad94f0>
```

```
In [90]: y_pred = classifier.predict(X_test)
y_pred = (y_pred > 0.5) #Predicting the result
```

```
In [97]: from sklearn.metrics import confusion_matrix,accuracy_score,classification_report
```

```
In [92]: cm = confusion_matrix(y_test,y_pred)
```

```
In [93]: cm
```

```
Out[93]: array([[2328,    72],  
   [ 425,  175]], dtype=int64)
```

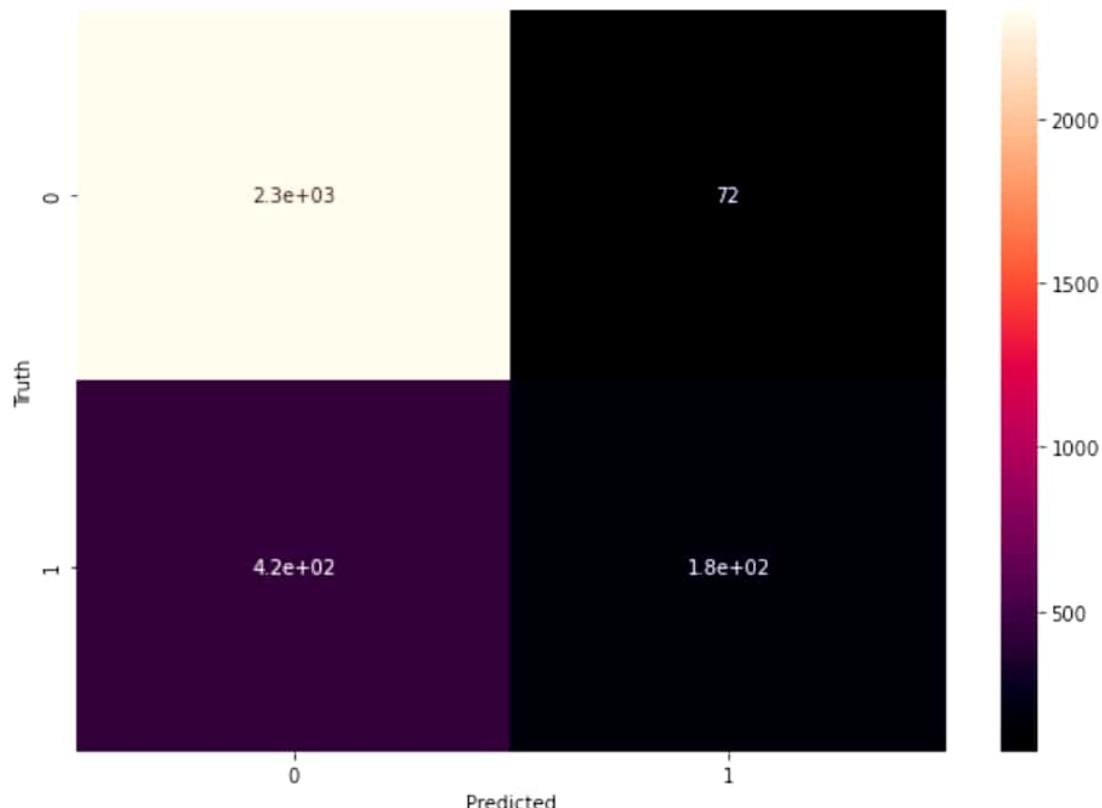
```
In [94]: accuracy = accuracy_score(y_test,y_pred)
```

```
In [95]: accuracy
```

```
Out[95]: 0.8343333333333334
```

```
In [98]: plt.figure(figsize = (10,7))  
sns.heatmap(cm,annot = True)  
plt.xlabel('Predicted')  
plt.ylabel('Truth')
```

```
Out[98]: Text(69.0, 0.5, 'Truth')
```



```
In [100... print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.85	0.97	0.90	2400
1	0.71	0.29	0.41	600
accuracy			0.83	3000
macro avg	0.78	0.63	0.66	3000
weighted avg	0.82	0.83	0.81	3000

```
In [ ]:
```

Implement K-Means clustering/ hierarchical clustering on sales_data_sample.csv dataset. Determine the number of clusters using the elbow method.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
#Importing the required Libraries.
```

```
In [2]: from sklearn.cluster import KMeans, k_means #For clustering
from sklearn.decomposition import PCA #Linear Dimensionality reduction.
```

```
In [3]: df = pd.read_csv("sales_data_sample.csv") #Loading the dataset.
```

Preprocessing

```
In [4]: df.head()
```

```
Out[4]: ORDERNUMBER QUANTITYORDERED PRICEEACH ORDERLINENUMBER SALES ORDERDATE S1
0 10107 30 95.70 2 2871.00 2/24/2003 0:00 Sh
1 10121 34 81.35 5 2765.90 5/7/2003 0:00 Sh
2 10134 41 94.74 2 3884.34 7/1/2003 0:00 Sh
3 10145 45 83.26 6 3746.70 8/25/2003 0:00 Sh
4 10159 49 100.00 14 5205.27 10/10/2003 0:00 Sh
```

5 rows × 25 columns

◀ ▶

```
In [5]: df.shape
```

```
Out[5]: (2823, 25)
```

```
In [6]: df.describe()
```

```
Out[6]: ORDERNUMBER QUANTITYORDERED PRICEEACH ORDERLINENUMBER SALES C
count 2823.000000 2823.000000 2823.000000 2823.000000 2823.000000 2823.000000
```

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	C
mean	10258.725115	35.092809	83.658544	6.466171	3553.889072	2.1
std	92.085478	9.741443	20.174277	4.225841	1841.865106	1.2
min	10100.000000	6.000000	26.880000	1.000000	482.130000	1.0
25%	10180.000000	27.000000	68.860000	3.000000	2203.430000	2.0
50%	10262.000000	35.000000	95.700000	6.000000	3184.800000	3.0
75%	10333.500000	43.000000	100.000000	9.000000	4508.000000	4.0
max	10425.000000	97.000000	100.000000	18.000000	14082.800000	4.0

◀ ▶

In [7]:

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   ORDERNUMBER      2823 non-null   int64  
 1   QUANTITYORDERED 2823 non-null   int64  
 2   PRICEEACH        2823 non-null   float64 
 3   ORDERLINENUMBER 2823 non-null   int64  
 4   SALES            2823 non-null   float64 
 5   ORDERDATE        2823 non-null   object  
 6   STATUS            2823 non-null   object  
 7   QTR_ID           2823 non-null   int64  
 8   MONTH_ID         2823 non-null   int64  
 9   YEAR_ID          2823 non-null   int64  
 10  PRODUCTLINE      2823 non-null   object  
 11  MSRP              2823 non-null   int64  
 12  PRODUCTCODE      2823 non-null   object  
 13  CUSTOMERNAME     2823 non-null   object  
 14  PHONE             2823 non-null   object  
 15  ADDRESSLINE1     2823 non-null   object  
 16  ADDRESSLINE2     302 non-null    object  
 17  CITY              2823 non-null   object  
 18  STATE             1337 non-null   object  
 19  POSTALCODE        2747 non-null   object  
 20  COUNTRY           2823 non-null   object  
 21  TERRITORY         1749 non-null   object  
 22  CONTACTLASTNAME   2823 non-null   object  
 23  CONTACTFIRSTNAME  2823 non-null   object  
 24  DEALSIZE          2823 non-null   object  
dtypes: float64(2), int64(7), object(16)
memory usage: 551.5+ KB
```

In [8]:

df.isnull().sum()

Out[8]:

ORDERNUMBER	0
QUANTITYORDERED	0
PRICEEACH	0
ORDERLINENUMBER	0
SALES	0
ORDERDATE	0
STATUS	0
QTR_ID	0
MONTH_ID	0
YEAR_ID	0
PRODUCTLINE	0
MSRP	0

```
PRODUCTCODE          0
CUSTOMERNAME        0
PHONE               0
ADDRESSLINE1         0
ADDRESSLINE2         2521
CITY                0
STATE               1486
POSTALCODE          76
COUNTRY             0
TERRITORY           1074
CONTACTLASTNAME    0
CONTACTFIRSTNAME   0
DEALSIZE            0
dtype: int64
```

```
In [9]: df.dtypes
```

```
Out[9]: ORDERNUMBER      int64
QUANTITYORDERED     int64
PRICEEACH           float64
ORDERLINENUMBER     int64
SALES              float64
ORDERDATE           object
STATUS              object
QTR_ID              int64
MONTH_ID             int64
YEAR_ID              int64
PRODUCTLINE         object
MSRP                int64
PRODUCTCODE          object
CUSTOMERNAME         object
PHONE               object
ADDRESSLINE1         object
ADDRESSLINE2         object
CITY                object
STATE               object
POSTALCODE           object
COUNTRY              object
TERRITORY            object
CONTACTLASTNAME     object
CONTACTFIRSTNAME    object
DEALSIZE             object
dtype: object
```

```
In [10]: df_drop  = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS','POSTALCODE', 'CITY', 'TERRITORY']
df = df.drop(df_drop, axis=1) #Dropping the categorical uneccessary columns along with it
```

```
In [11]: df.isnull().sum()
```

```
Out[11]: QUANTITYORDERED    0
PRICEEACH           0
ORDERLINENUMBER     0
SALES              0
ORDERDATE           0
QTR_ID              0
MONTH_ID             0
YEAR_ID              0
PRODUCTLINE         0
MSRP                0
PRODUCTCODE          0
COUNTRY             0
DEALSIZE            0
dtype: int64
```

```
In [12]:
```

```

df.dtypes

Out[12]: QUANTITYORDERED      int64
PRICEEACH          float64
ORDERLINENUMBER    int64
SALES             float64
ORDERDATE          object
QTR_ID            int64
MONTH_ID           int64
YEAR_ID            int64
PRODUCTLINE        object
MSRP              int64
PRODUCTCODE        object
COUNTRY            object
DEALSIZE           object
dtype: object

In [13]: # Checking the categorical columns.

In [14]: df['COUNTRY'].unique()

Out[14]: array(['USA', 'France', 'Norway', 'Australia', 'Finland', 'Austria', 'UK',
   'Spain', 'Sweden', 'Singapore', 'Canada', 'Japan', 'Italy',
   'Denmark', 'Belgium', 'Philippines', 'Germany', 'Switzerland',
   'Ireland'], dtype=object)

In [15]: df['PRODUCTLINE'].unique()

Out[15]: array(['Motorcycles', 'Classic Cars', 'Trucks and Buses', 'Vintage Cars',
   'Planes', 'Ships', 'Trains'], dtype=object)

In [16]: df['DEALSIZE'].unique()

Out[16]: array(['Small', 'Medium', 'Large'], dtype=object)

In [17]: productline = pd.get_dummies(df['PRODUCTLINE']) #Converting the categorical columns.
Dealsize = pd.get_dummies(df['DEALSIZE'])

In [18]: df = pd.concat([df,productline,Dealsize], axis = 1)

In [19]: df_drop  = ['COUNTRY','PRODUCTLINE','DEALSIZE'] #Dropping Country too as there are a
df = df.drop(df_drop, axis=1)

In [20]: df['PRODUCTCODE'] = pd.Categorical(df['PRODUCTCODE']).codes #Converting the datatype

In [21]: df.drop('ORDERDATE', axis=1, inplace=True) #Dropping the Orderdate as Month is alrea

In [22]: df.dtypes #All the datatypes are converted into numeric

Out[22]: QUANTITYORDERED      int64
PRICEEACH          float64
ORDERLINENUMBER    int64
SALES             float64

```

```

QTR_ID           int64
MONTH_ID         int64
YEAR_ID          int64
MSRP             int64
PRODUCTCODE      int8
Classic Cars    uint8
Motorcycles      uint8
Planes           uint8
Ships             uint8
Trains            uint8
Trucks and Buses uint8
Vintage Cars     uint8
Large              uint8
Medium             uint8
Small              uint8
dtype: object

```

Plotting the Elbow Plot to determine the number of clusters.

In [23]:

```

distortions = [] # Within Cluster Sum of Squares from the centroid
K = range(1,10)
for k in K:
    kmeanModel = KMeans(n_clusters=k)
    kmeanModel.fit(df)
    distortions.append(kmeanModel.inertia_) #Appending the inertia to the Distortion

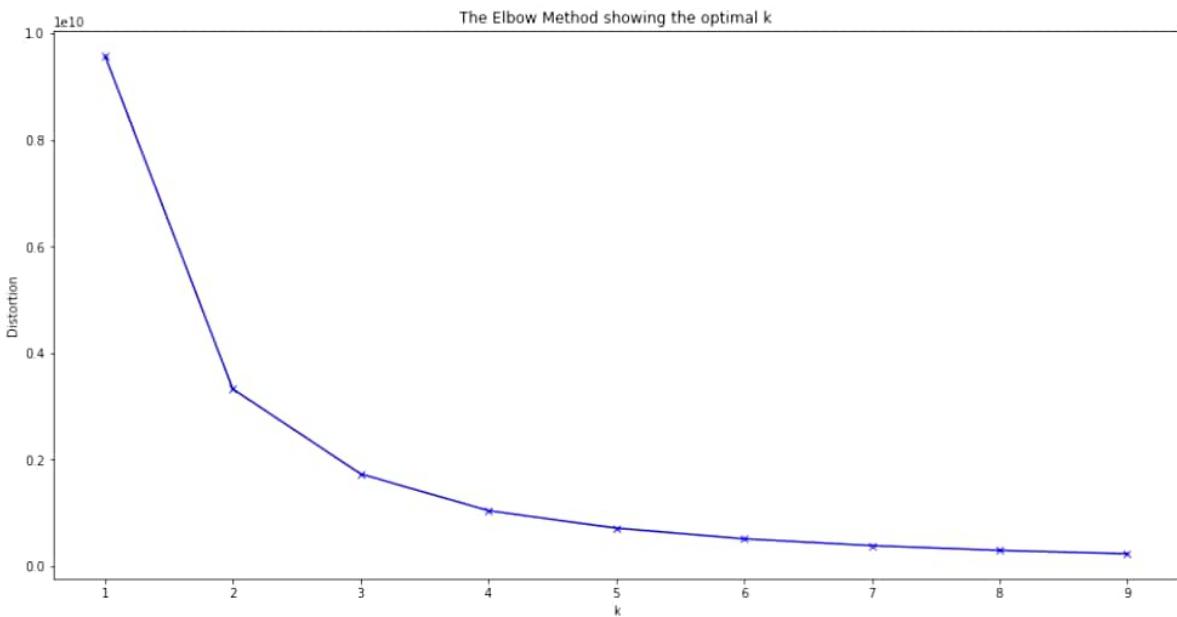
```

In [24]:

```

plt.figure(figsize=(16,8))
plt.plot(K, distortions, 'bx-')
plt.xlabel('k')
plt.ylabel('Distortion')
plt.title('The Elbow Method showing the optimal k')
plt.show()

```



As the number of k increases Inertia decreases.

Observations: A Elbow can be observed at 3 and after that the curve decreases gradually.

```
In [25]: X_train = df.values #Returns a numpy array.
```

```
In [26]: X_train.shape
```

```
Out[26]: (2823, 19)
```

```
In [27]: model = KMeans(n_clusters=3,random_state=2) #Number of cluster = 3
model = model.fit(X_train) #Fitting the values to create a model.
predictions = model.predict(X_train) #Predicting the cluster values (0,1,or 2)
```

```
In [28]: unique,counts = np.unique(predictions,return_counts=True)
```

```
In [29]: counts = counts.reshape(1,3)
```

```
In [30]: counts_df = pd.DataFrame(counts,columns=['Cluster1','Cluster2','Cluster3'])
```

```
In [31]: counts_df.head()
```

```
Out[31]:
```

	Cluster1	Cluster2	Cluster3
0	1367	373	1083

Visualization

```
In [32]: pca = PCA(n_components=2) #Converting all the features into 2 columns to make it eas
```

```
In [33]: reduced_X = pd.DataFrame(pca.fit_transform(X_train),columns=['PCA1','PCA2']) #Creati
```

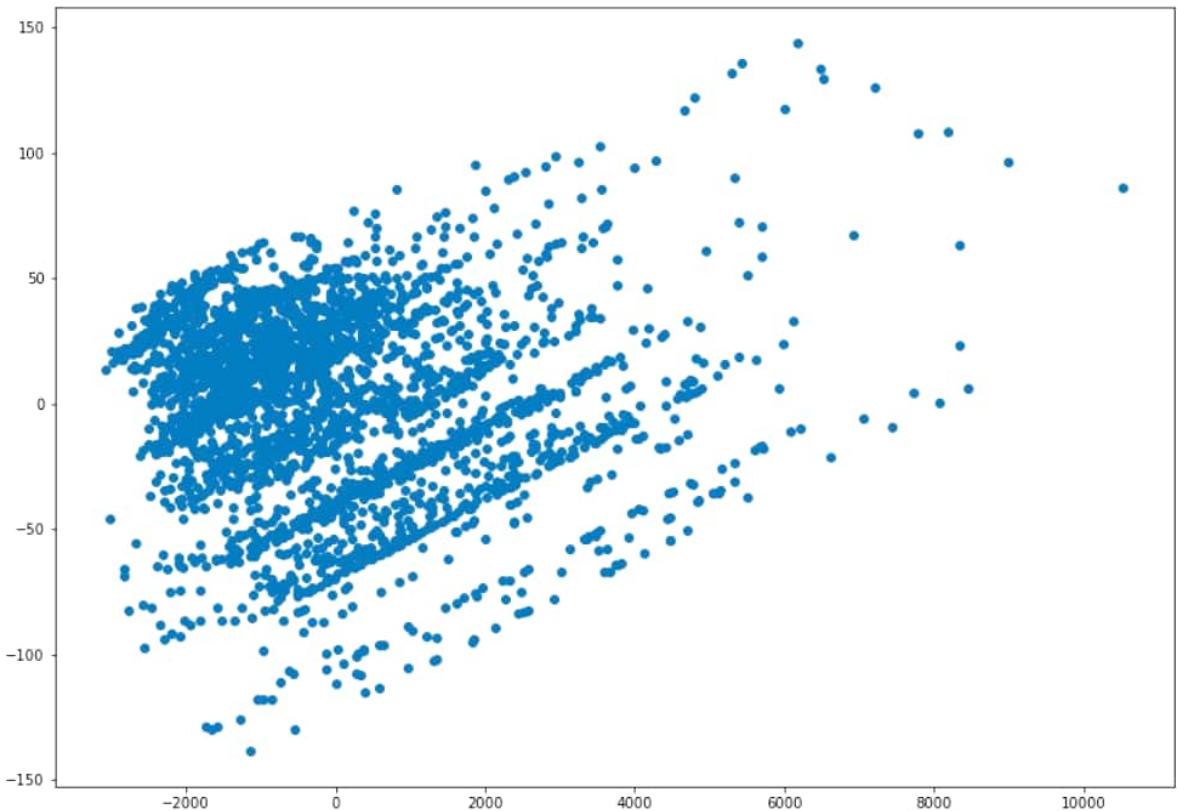
```
In [34]: reduced_X.head()
```

```
Out[34]:
```

	PCA1	PCA2
0	-682.488323	-42.819535
1	-787.665502	-41.694991
2	330.732170	-26.481208
3	193.040232	-26.285766
4	1651.532874	-6.891196

```
In [35]: #Plotting the normal Scatter Plot
plt.figure(figsize=(14,10))
plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
```

```
Out[35]: <matplotlib.collections.PathCollection at 0x2b3b54c6f70>
```



```
In [36]: model.cluster_centers_ #Finding the centriods. (3 Centriods in total. Each Array con
```

```
Out[36]: array([[ 3.08302853e+01,  7.00755230e+01,  6.67300658e+00,
   2.12409474e+03,  2.71762985e+00,  7.09509876e+00,
   2.00381127e+03,  7.84784199e+01,  6.24871982e+01,
   2.64813460e-01,  1.21433797e-01,  1.29480614e-01,
   1.00219459e-01,  3.87710315e-02,  9.21726408e-02,
   2.53108998e-01,  6.93889390e-18,  6.21799561e-02,
   9.37820044e-01],
 [ 4.45871314e+01,  9.98931099e+01,  5.75603217e+00,
   7.09596863e+03,  2.71045576e+00,  7.06434316e+00,
   2.00389008e+03,  1.45823056e+02,  3.14959786e+01,
   5.33512064e-01,  1.07238606e-01,  7.23860590e-02,
   2.14477212e-02,  1.07238606e-02,  1.31367292e-01,
   1.23324397e-01,  4.20911528e-01,  5.79088472e-01,
   5.55111512e-17],
 [ 3.72031394e+01,  9.52120960e+01,  6.44967682e+00,
   4.13868425e+03,  2.72022161e+00,  7.09879963e+00,
   2.00379409e+03,  1.13248384e+02,  5.04469067e+01,
   3.74884580e-01,  1.15420129e-01,  9.41828255e-02,
   8.21791320e-02,  1.84672207e-02,  1.16343490e-01,
   1.98522622e-01,  2.08166817e-17,  1.00000000e+00,
   -6.66133815e-16]])
```

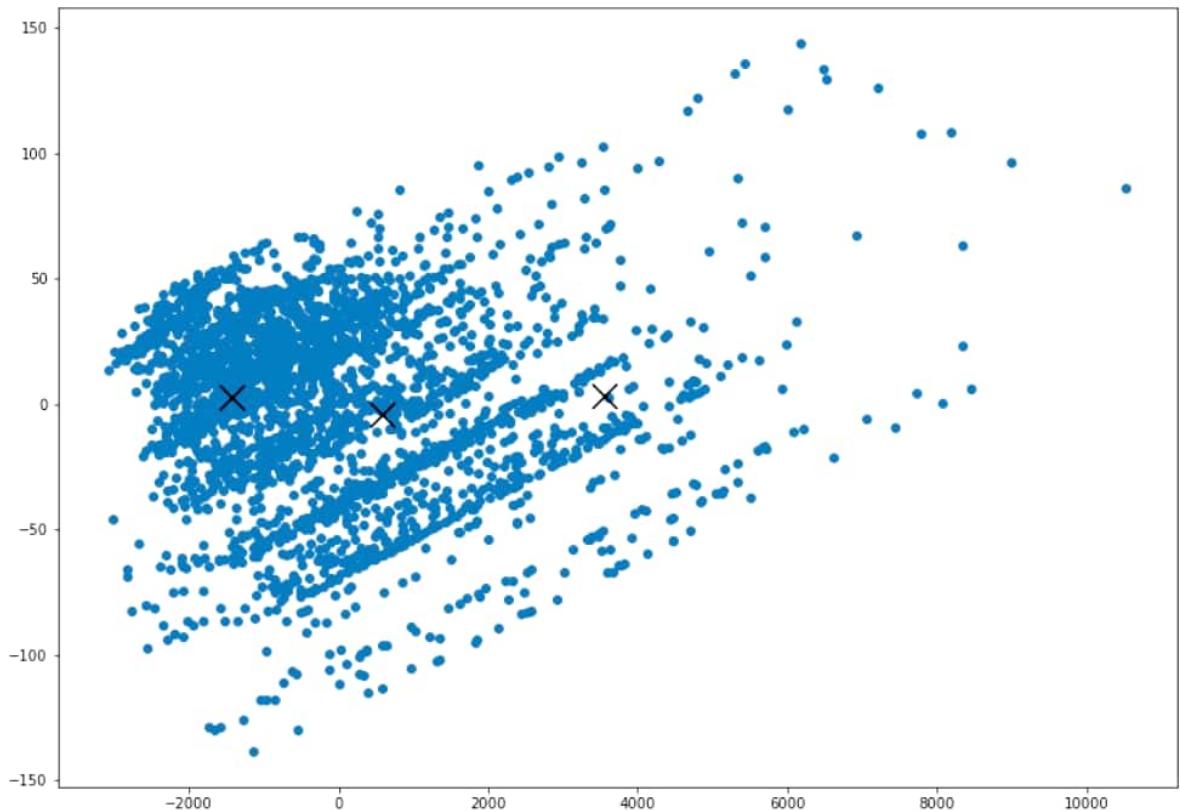
```
In [37]: reduced_centers = pca.transform(model.cluster_centers_) #Transforming the centroids
```

```
In [38]: reduced_centers
```

```
Out[38]: array([[-1.43005891e+03,  2.60041009e+00],
 [ 3.54247180e+03,  3.15185487e+00],
 [ 5.84994044e+02, -4.36786931e+00]])
```

```
In [39]: plt.figure(figsize=(14,10))
plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker='x',s=300)
```

```
Out[39]: <matplotlib.collections.PathCollection at 0x2b3b555b130>
```



```
In [40]: reduced_X['Clusters'] = predictions #Adding the Clusters to the reduced dataframe.
```

```
In [41]: reduced_X.head()
```

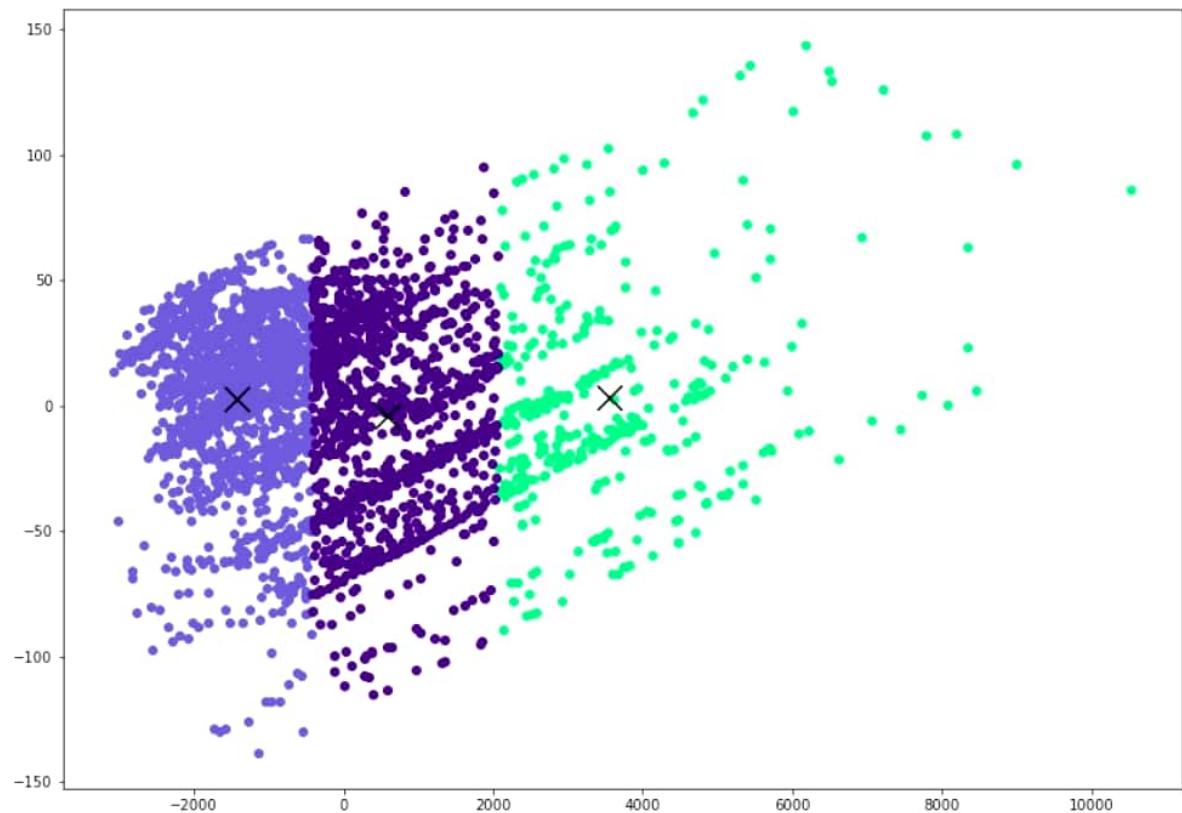
```
Out[41]:
```

	PCA1	PCA2	Clusters
0	-682.488323	-42.819535	0
1	-787.665502	-41.694991	0
2	330.732170	-26.481208	2
3	193.040232	-26.285766	2
4	1651.532874	-6.891196	2

```
In [42]: #Plotting the clusters
plt.figure(figsize=(14,10))
# taking the cluster number and first column
# taking the cluster number and second column
# taking the cluster number and third column
plt.scatter(reduced_X[reduced_X['Clusters'] == 0].loc[:, 'PCA1'], reduced_X[reduced_X['Clusters'] == 0].loc[:, 'PCA2'])
plt.scatter(reduced_X[reduced_X['Clusters'] == 1].loc[:, 'PCA1'], reduced_X[reduced_X['Clusters'] == 1].loc[:, 'PCA2'])
plt.scatter(reduced_X[reduced_X['Clusters'] == 2].loc[:, 'PCA1'], reduced_X[reduced_X['Clusters'] == 2].loc[:, 'PCA2'])

plt.scatter(reduced_centers[:,0], reduced_centers[:,1], color='black', marker='x', s=300)
```

```
Out[42]: <matplotlib.collections.PathCollection at 0x2b3b5778af0>
```



In []:

In []: