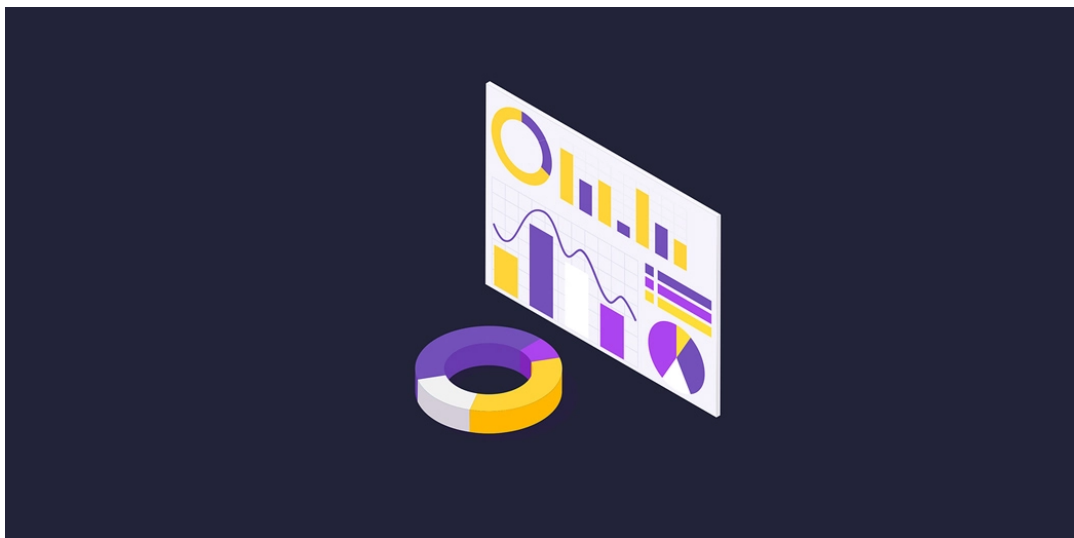


Exploratory Data Analysis(EDA) in Python!

Introduction

Exploratory data analysis is one of the best practices used in data science today. While starting a career in Data Science, people generally don't know the difference between Data analysis and exploratory data analysis. There is not a very big difference between the two, but both have different purposes.



Exploratory Data Analysis(EDA): Exploratory data analysis is a complement to [inferential statistics](#), which tends to be fairly rigid with rules and formulas. At an advanced level, EDA involves looking at and describing the data set from different angles and then summarizing it.

Data Analysis: Data Analysis is the statistics and probability to figure out trends in the data set. It is used to show historical data by using some analytics tools. It helps in drilling down the information, to transform metrics, facts, and figures into initiatives for improvement.

Exploratory Data Analysis(EDA)

We will explore a Data set and perform the exploratory data analysis in python. You can refer to our [python course online](#) to get on board with python.

The major topics to be covered are below:

- Handle Missing value
- Removing duplicates
- Outlier Treatment
- Normalizing and Scaling(Numerical Variables)
- Encoding Categorical variables(Dummy Variables)
- Bivariate Analysis

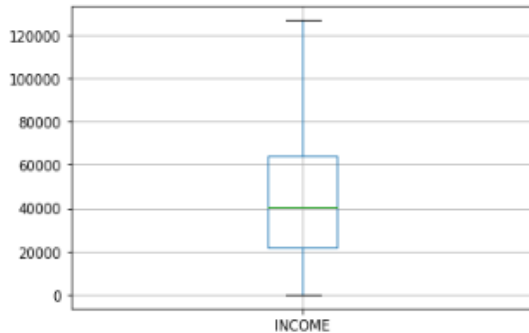
Importing Libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

Loading the data set

We will be loading the EDA cars excel file using pandas. For this, we will be using read_excel file.

```
car_df.boxplot(column=['INCOME'])
plt.show()
```



Box-plot after removing outliers

Basic Data Exploration

In this step, we will perform the below operations to check what the data set comprises of. We will check the below things:

- head of the dataset
- the shape of the dataset
- info of the dataset
- summary of the dataset

1. The head function will tell you the top records in the data set. By default, python shows you only the top 5 records.
2. The shape attribute tells us a number of observations and variables we have in the data set. It is used to check the dimension of data. The cars data set has 303 observations and 13 variables in the data set.

```
car_df.shape
(303, 13)
```

3. info() is used to check the Information about the data and the datatypes of each respective attribute.

```
car_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 13 columns):
INDEX            303 non-null int64
INCOME           260 non-null float64
MARITAL STATUS   275 non-null object
SEX              297 non-null object
EDUCATION        259 non-null object
JOB              257 non-null object
TRAVEL TIME      262 non-null float64
USE              250 non-null object
MILES CLOCKED    278 non-null float64
CAR TYPE         293 non-null object
CAR AGE          283 non-null float64
CITY             297 non-null object
POSTAL CODE      300 non-null float64
dtypes: float64(5), int64(1), object(7)
memory usage: 30.9+ KB
```

Looking at the data in the head function and in info, we know that the variable Income and travel time are of float data type instead of the object. So we will convert it into the float. Also, there are some invalid values like @@ and '*' in the data which we will be treating as missing values.

```
car_df['INCOME'] = car_df['INCOME'].replace(to_replace = '@@', value = np.nan)
car_df['INCOME'] = car_df['INCOME'].astype(float)

car_df['TRAVEL TIME'] = car_df['TRAVEL TIME'].replace(to_replace = '*****', value = np.nan)

car_df['TRAVEL TIME'] = car_df['TRAVEL TIME'].astype(float)

car_df['MILES CLOCKED'] = car_df['MILES CLOCKED'].replace(to_replace = 'Na', value = np.nan)
car_df['MILES CLOCKED'] = car_df['MILES CLOCKED'].replace(to_replace = '*****', value = np.nan)
car_df['MILES CLOCKED'] = car_df['MILES CLOCKED'].astype(float)

car_df.replace(to_replace = '*****', value = np.nan, inplace=True)
```

```
car_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 13 columns):
INDEX            303 non-null int64
INCOME           260 non-null float64
MARITAL STATUS   275 non-null object
SEX              297 non-null object
EDUCATION        259 non-null object
JOB              257 non-null object
TRAVEL TIME      262 non-null float64
USE              250 non-null object
MILES CLOCKED    278 non-null float64
CAR TYPE         293 non-null object
CAR AGE          283 non-null float64
CITY             297 non-null object
POSTAL CODE      300 non-null float64
dtypes: float64(5), int64(1), object(7)
memory usage: 30.9+ KB
```

- The described method will help to see how data has been spread for numerical values. We can clearly see the minimum value, mean values, different percentile values, and maximum values.

```
car_df.describe()
```

	INDEX	INCOME	TRAVEL TIME	MILES CLOCKED	CAR AGE	POSTAL CODE
count	303.000000	260.000000	262.000000	278.000000	283.000000	300.000000
mean	139.640264	50025.162170	34.282098	13591.978417	6.265018	50712.198667
std	85.178422	41188.807914	14.910178	7167.328655	5.111218	24141.029290
min	1.000000	0.000000	5.000000	1500.000000	1.000000	11435.000000
25%	62.500000	20452.885022	24.449874	7900.000000	1.000000	42420.000000
50%	138.000000	44571.590870	33.564757	12065.000000	6.000000	47150.000000
75%	213.500000	66485.761387	43.907339	18240.000000	10.000000	61701.000000
max	289.000000	204667.589700	83.617643	38000.000000	20.000000	80049.000000

Handling missing value

```
# Check for missing value in any column
car_df.isnull().sum()

INDEX          0
INCOME         43
MARITAL STATUS 28
SEX            6
EDUCATION      44
JOB            46
TRAVEL TIME    41
USE            53
MILES CLOCKED  25
CAR TYPE       10
CAR AGE        20
CITY           6
POSTAL CODE    3
dtype: int64
```

We can see that we have various missing values in the respective columns. There are various ways of treating your missing values in the data set. And which technique to use when is actually dependent on the type of data you are dealing with.

- Drop the missing values: In this case, we drop the missing values from those variables. In case there are very few missing values you can drop those values.
- Impute with mean value: For the numerical column, you can replace the missing values with mean values. Before replacing with mean value, it is advisable to check that the variable shouldn't have extreme values .i.e. outliers.
- Impute with median value: For the numerical column, you can also replace the missing values with median values. In case you have extreme values such as outliers it is advisable to use the median approach.
- Impute with mode value: For the categorical column, you can replace the missing values with mode values i.e the frequent ones.

In this exercise, we will replace the numerical columns with median values and for categorical columns, we will drop the missing values.

```
# Replacing NULL values in Numerical Columns using Median
median1=car_df["INCOME"].median()
median2=car_df["TRAVEL TIME"].median()
median3=car_df["MILES CLOCKED"].median()
median4=car_df["CAR AGE"].median()
median5=car_df["POSTAL CODE"].median()

car_df["INCOME"].replace(np.nan,median1,inplace=True)
car_df["TRAVEL TIME"].replace(np.nan,median2,inplace=True)
car_df["MILES CLOCKED"].replace(np.nan,median3,inplace=True)
car_df["CAR AGE"].replace(np.nan,median4,inplace=True)
car_df["POSTAL CODE"].replace(np.nan,median5,inplace=True)
```

```
# Replacing NULL values in catogrical Columns using mode
mode1=car_df["SEX"].mode().values[0]
mode2=car_df["MARITAL STATUS"].mode().values[0]
mode3=car_df["EDUCATION"].mode().values[0]
mode4=car_df["JOB"].mode().values[0]
mode5=car_df["USE"].mode().values[0]
mode6=car_df['CITY'].mode().values[0]
mode7=car_df["CAR TYPE"].mode().values[0]

car_df["SEX"]=car_df["SEX"].replace(np.nan,mode1)
car_df["MARITAL STATUS"]= car_df["MARITAL STATUS"].replace(np.nan,mode2)
car_df["EDUCATION"]=car_df["EDUCATION"].replace(np.nan,mode3)
car_df["JOB"]=car_df["JOB"].replace(np.nan,mode4)
car_df["USE"]=car_df["USE"].replace(np.nan,mode5)
car_df["CAR TYPE"]=car_df["CAR TYPE"].replace(np.nan,mode7)
car_df['CITY']=car_df['CITY'].replace(np.nan,mode6)
```

```
car_df.isnull().sum()
```

```
INDEX      0
INCOME      0
MARITAL STATUS  0
SEX         0
EDUCATION   0
JOB         0
TRAVEL TIME  0
USE         0
MILES CLOCKED  0
CAR TYPE    0
CAR AGE     0
CITY        0
POSTAL CODE  0
dtype: int64
```

Handling Duplicate records

```
# Check for duplicate data
duplicate = car_df.duplicated()
print(duplicate.sum())
car_df[duplicate]
```

14

	INDEX	INCOME	MARITAL STATUS	SEX	EDUCATION	JOB	TRAVEL TIME	USE	MILES CLOCKED	CAR TYPE	CAR AGE	CITY	POSTAL CODE
69	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
70	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
71	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
72	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
73	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
74	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
75	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
76	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
77	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
78	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
79	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
80	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
81	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0
82	29	64013.81632	Yes	M	High School	Blue Collar	32.717234	Commercial	7900.0	Pickup	5.0	Los Angeles	90049.0

Since we have 14 duplicate records in the data, we will remove this from the data set so that we get only distinct records. Post removing the duplicate, we will check whether the duplicates have been removed from the data set or not.

```
#write the code to drop the duplicates.
car_df.drop_duplicates(inplace=True)
```

```
# Check for duplicate data again

dpl = car_df.duplicated()
dpl.sum()
```

0

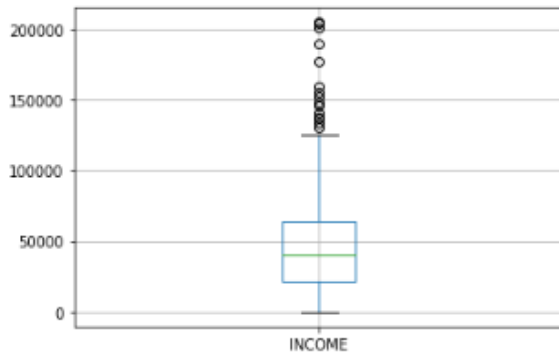
Handling Outlier

Outliers, being the most extreme observations, may include the sample maximum or sample minimum, or both, depending on whether they are extremely high or low. However, the sample maximum and minimum are not always outliers because they may not be unusually far from other observations.

We Generally identify outliers with the help of boxplot, so here box plot shows some of the data points outside the range of the data.

```
car_df.boxplot(column=["INCOME"])
plt.show
```

```
<function matplotlib.pyplot.show(*args, **kw)>
```



Box-plot before removing outliers

Looking at the box plot, it seems that the variables INCOME, have outlier present in the variables. These outliers value needs to be teated and there are several ways of treating them:

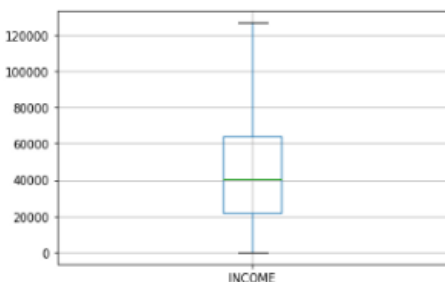
- Drop the outlier value
- Replace the outlier value using the IQR

```
#def remove_outlier(col):
#create a user definded function called remove_outlier for getting the threshold value from IQR.
def remove_outlier(col):
    sorted(col)
    Q1,Q3=col.quantile([0.25,0.75])
    IQR=Q3-Q1
    lower_range= Q1-(1.5 * IQR)
    upper_range= Q3+(1.5 * IQR)
    return lower_range, upper_range
```

```
lowincome,uppincome=remove_outlier(car_df['INCOME'])
car_df['INCOME']=np.where(car_df['INCOME']>uppincome,uppincome,car_df['INCOME'])
car_df['INCOME']=np.where(car_df['INCOME']<lowincome,lowincome,car_df['INCOME'])
```

#Boxplot After removing outlier

```
car_df.boxplot(column=['INCOME'])
plt.show()
```



Bivariate Analysis

When we talk about bivariate analysis, it means analyzing 2 variables. Since we know there are numerical and categorical variables, there is a way of analyzing these variables as shown below:

1. Numerical vs. Numerical

1. Scatterplot
2. Line plot
3. Heatmap for correlation
4. Joint plot

2. Categorical vs. Numerical

1. Bar chart
2. Violin plot
3. Categorical box plot
4. Swarm plot

3. Two Categorical Variables

1. Bar chart
2. Grouped bar chart
3. Point plot

If we need to find the correlation-

```
car_df.corr()
```

	INDEX	INCOME	TRAVEL TIME	MILES CLOCKED	CAR AGE	POSTAL CODE
INDEX	1.000000	-0.041232	0.022397	0.038455	-0.027206	-0.244129
INCOME	-0.041232	1.000000	0.061243	0.340653	0.263088	0.035855
TRAVEL TIME	0.022397	0.061243	1.000000	0.023395	0.143684	0.017591
MILES CLOCKED	0.038455	0.340653	0.023395	1.000000	0.130708	-0.113445
CAR AGE	-0.027206	0.263088	0.143684	0.130708	1.000000	-0.098372
POSTAL CODE	-0.244129	0.035855	0.017591	-0.113445	-0.098372	1.000000

Correlation between all the variables

Normalizing and Scaling

Often the variables of the data set are of different scales i.e. one variable is in millions and others in only 100. For e.g. in our data set Income is having values in thousands and age in just two digits. Since the data in these variables are of different scales, it is tough to compare these variables.

Feature scaling (also known as data normalization) is the method used to standardize the range of features of data. Since the range of values of data may vary widely, it becomes a necessary step in data preprocessing while using machine learning algorithms.

In this method, we convert variables with different scales of measurements into a single scale. StandardScaler normalizes the data using the formula $(x - \text{mean}) / \text{standard deviation}$. We will be doing this only for the numerical variables.

```
#Scales the data. Essentially returns the z-scores of every attribute
# z-score:z-score (also called a standard score) gives you an idea of how far from the mean a data point is.
#But more technically it's a measure of how many standard deviations below or above the population mean a raw score is
from sklearn.preprocessing import StandardScaler
std_scale = StandardScaler()
std_scale

StandardScaler(copy=True, with_mean=True, with_std=True)
```

```
car_df['INCOME'] = std_scale.fit_transform(car_df[['INCOME']])
car_df['TRAVEL TIME'] = std_scale.fit_transform(car_df[['TRAVEL TIME']])
car_df['CAR AGE'] = std_scale.fit_transform(car_df[['CAR AGE']])
car_df['POSTAL CODE'] = std_scale.fit_transform(car_df[['POSTAL CODE']])
car_df['MILES CLOCKED'] = std_scale.fit_transform(car_df[['MILES CLOCKED']])
```

```
car_df.head()
```

	INDEX	INCOME	MARITAL STATUS	SEX	EDUCATION	JOB	TRAVEL TIME	USE	MILES CLOCKED	CAR TYPE	CAR AGE	CITY	POSTAL CODE
0	1	2.344391	No	F	Bachelors	Blue Collar	0.832200	Commercial	0.535034	Sports Car	0.137267	Texas	-0.276006
1	2	0.136993	No	M	High School	Blue Collar	-0.991124	Private	0.754798	Minivan	-1.052842	Texas	-0.276006
2	3	0.498435	No	F	Bachelors	Clerical	-0.043693	Private	-0.136620	SUV	-1.052842	Texas	-0.276006
3	4	0.917492	No	F	High School	Lawyer	-1.366047	Private	0.662480	Sports Car	0.930674	Texas	-0.276006
4	5	2.391357	No	M	High School	Blue Collar	0.013415	Commercial	2.133234	Panel Truck	0.732322	Texas	-0.276006

ENCODING

One-Hot-Encoding is used to create dummy variables to replace the categories in a categorical variable into features of each category and represent it using 1 or 0 based on the presence or absence of the categorical value in the record.

This is required to do since the machine learning algorithms only work on the numerical data. That is why there is a need to convert the categorical column into a numerical one.

`get_dummies` is the method that creates a dummy variable for each categorical variable.

```
dummies=pd.get_dummies(car_df[["MARITAL STATUS", "SEX","EDUCATION","JOB","USE","CAR TYPE","CITY"]],
                           columns=["MARITAL STATUS", "SEX","EDUCATION","JOB","USE","CAR TYPE","CITY"],
                           prefix=["married", "sex", "Education", "Job", "Use", "cartype", "city"],drop_first=True).head()
```

```
dummies.head()
```

	married_Yes	sex_M	Education_High School	Education_Masters	Education_PhD	Job_Clerical	Job_Doctor	Job_Home Maker	Job_Lawyer	Job_Manager	...	city_Houston
0	0	0	0	0	0	0	0	0	0	0	...	0
1	0	1	1	0	0	0	0	0	0	0	...	0
2	0	0	0	0	0	1	0	0	0	0	...	0
3	0	0	1	0	0	0	0	0	1	0	...	0
4	0	1	1	0	0	0	0	0	0	0	...	0

5 rows x 28 columns

```
columns=["MARITAL STATUS", "SEX","EDUCATION","JOB","USE","CAR TYPE","CITY"]
car_df = pd.concat([car_df, dummies], axis=1)
# drop original column "fuel-type" from "df"
car_df.drop(columns, axis = 1, inplace=True)
```

```
car_df.head()
```

	INDEX	INCOME	TRAVEL TIME	MILES CLOCKED	CAR AGE	POSTAL CODE	married_Yes	sex_M	Education_High School	Education_Masters	...	city_Houston	city_Las Vegas	city_Los Angeles
0	1	2.344391	0.832200	0.535034	0.137267	-0.276006	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0
1	2	0.136993	-0.991124	0.754798	-1.052842	-0.276006	0.0	1.0	1.0	0.0	...	0.0	0.0	0.0
2	3	0.498435	-0.043693	-0.136620	-1.052842	-0.276006	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0
3	4	0.917492	-1.366047	0.662480	0.930674	-0.276006	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0
4	5	2.391357	0.013415	2.133234	0.732322	-0.276006	0.0	1.0	1.0	0.0	...	0.0	0.0	0.0

5 rows x 34 columns

About the Author



Ritika Singh – Data Scientist

I am a Data scientist by profession and a Blogger by passion. I have been working on machine learning projects for more than 2 years. Here you will find articles on “Machine Learning, Statistics, Deep Learning, NLP and Artificial Intelligence”.

