# EV Registration Analysis: Assignment

Data Science: Python Project: DATA Cleaning

Name: Chumani Mazwi

Email: cmazwi5@gmail.com

Institution: Deviare| MICT SETA| Belong|

Clean the data by removing incorrect, duplicate or missing data.

```
3   BONNEVILLE POWER ADMINISTRATION||PUD NO 1 OF C...     5.301104e+10
4                            PUGET SOUND ENERGY INC        5.303508e+10

[5]: # 1. Handling Missing Values

    df = df.dropna()
    df = df.fillna(0)

[6]: df.isnull().sum()

[6]: Identifier                                           0
     City                                                 0
     Postal Code                                          0
     Model Year                                           0
     Make                                                 0
     Model                                                0
     Electric Vehicle Type                                0
     Clean Alternative Fuel Vehicle (CAFV) Eligibility    0
     Electric Range                                       0
     Base MSRP                                            0
     Legislative District                                 0
     Vehicle ID                                           0
     Vehicle Location                                     0
     Electric Utility                                     0
     2020 Census Tract                                    0
     dtype: int64

[ ]:
```



```
     dtype: int64

[7]: #Size of original dataset
     print(df.shape)

     (134474, 15)

[8]: #Dropping the missing rows.
     df_dropped = df.dropna(how = 'any')

[11]: # 2. Removing Duplicates
      df = df.drop_duplicates()

[13]: df=pd.read_csv("EV_Registration_Dataset.csv")

[14]: #checking the duplicates
      df.duplicated().sum()

[14]: 0

[ ]:
```
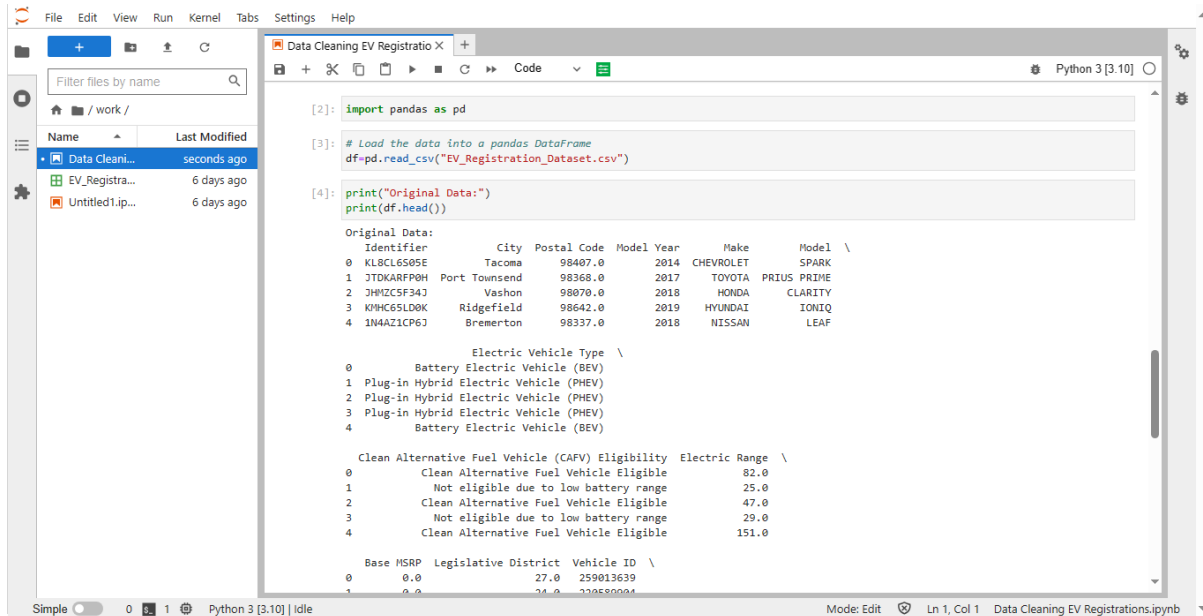
Clean the data by removing incorrect, duplicate or missing data:

Step 1:

Importing the necessary libraries and loading the data set:

```
[2]: import pandas as pd

[3]: # Load the data into a pandas DataFrame
     df=pd.read_csv("EV_Registration_Dataset.csv")

[4]: print("Original Data:")
     print(df.head())

     Original Data:
         Identifier          City  Postal Code Model Year      Make       Model  \
     0  KL8CL6S05E        Tacoma       98407.0       2014  CHEVROLET       SPARK
     1  JTDKARFP0H  Port Townsend     98368.0       2017     TOYOTA  PRIUS PRIME
     2  JHMZC5F34J        Vashon       98070.0       2018      HONDA     CLARITY
     3  KMHC65LD0K     Ridgefield      98642.0       2019    HYUNDAI       IONIQ
     4  1N4AZ1CP6J      Bremerton      98337.0       2018     NISSAN        LEAF

                            Electric Vehicle Type  \
     0            Battery Electric Vehicle (BEV)
     1  Plug-in Hybrid Electric Vehicle (PHEV)
     2  Plug-in Hybrid Electric Vehicle (PHEV)
     3  Plug-in Hybrid Electric Vehicle (PHEV)
     4            Battery Electric Vehicle (BEV)

       Clean Alternative Fuel Vehicle (CAFV) Eligibility  Electric Range  \
     0           Clean Alternative Fuel Vehicle Eligible            82.0
     1               Not eligible due to low battery range           25.0
     2           Clean Alternative Fuel Vehicle Eligible            47.0
     3               Not eligible due to low battery range           29.0
     4           Clean Alternative Fuel Vehicle Eligible           151.0

        Base MSRP  Legislative District  Vehicle ID  \
     0        0.0                  27.0   259013639
     1        0.0                  24.0   229580904
```

Step 2: Data Cleaning Steps: 1 { # Handling Missing Values}

Dropping rows with missing values| Alternatively, fill missing values with a specific value e.g (0):

```
[5]: #Dropping rows with missing Values:
     df = df.dropna()

[7]: # df = df.fillna(0)
     df=df.fillna(0)
```

Step 3: Removing Duplicates and Printing the Cleaned Data

```
[8]: #Removing Duplicates
     df = df.drop_duplicates()

[9]: print("\nCleaned Data:")

     Cleaned Data:

[10]: print(df.head())

         Identifier          City  Postal Code Model Year      Make       Model  \
     0  KL8CL6S05E        Tacoma       98407.0       2014  CHEVROLET       SPARK
     1  JTDKARFP0H  Port Townsend     98368.0       2017     TOYOTA  PRIUS PRIME
     2  JHMZC5F34J        Vashon       98070.0       2018      HONDA     CLARITY
     3  KMHC65LD0K     Ridgefield      98642.0       2019    HYUNDAI       IONIQ
```
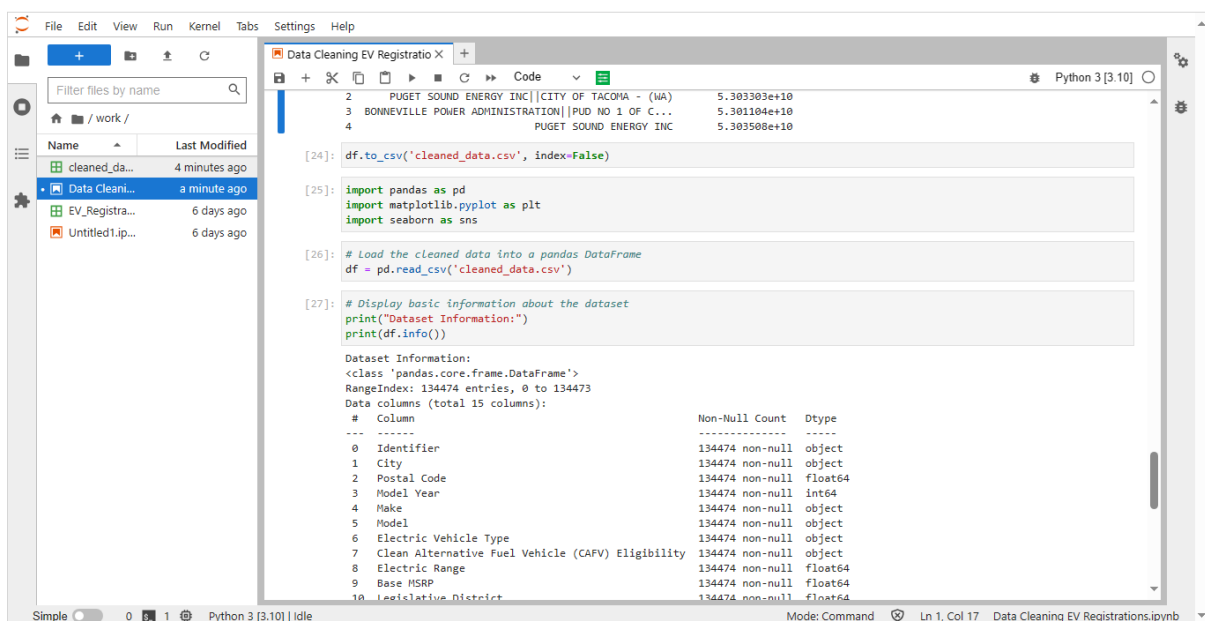
Analyse the data and understand the nature of the trends present in the data:

The first step is to import the necessary libraries, load the data into the pandas Data frame and display the basic information about the dataset:

The next step involves applying descriptive statistics about the dataset:



The next step involves using the correlation matrix to interpret the dataset:



The next step involves using the pair plot to interpret the dataset:

```
[30]: # Pairplot
      sns.pairplot(df)
      plt.title('Pairplot')
      plt.show()
```