# Data Collection and Preprocessing Phase

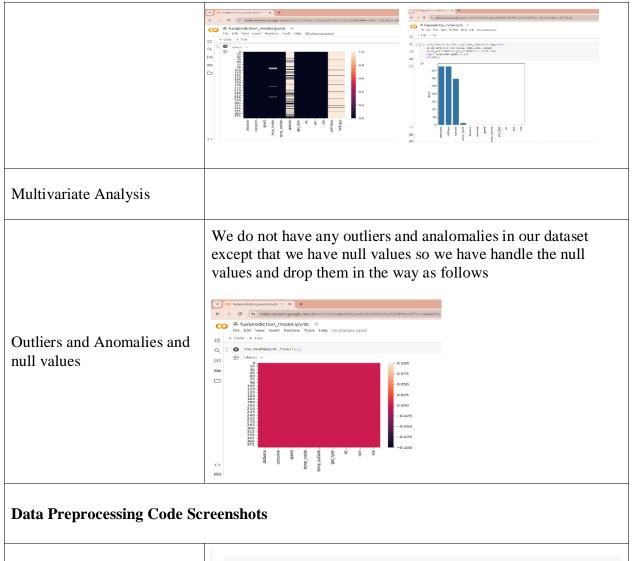| | |
|---|---|
| Date | 9 July 2024 |
| Team ID | 739659 |
| Project Title | Trip-Based Modelling of Fuel Consumption in Modern Fleet Vehicles Using Machine Learning |
| Maximum Marks | 6 Marks |

## Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview | 388 rows x 12 columns, dtypes: float64(4),int64(5),object(3)<br> |
| Univariate Analysis | Exploration of individual of accuracy_score,mean_squared_error,r2_score,mean_absolute_error |
| Bivariate Analysis | Relationships between two variables (correlation, scatter plots) . |

| | |
|---|---|
| Multivariate Analysis | |
| Outliers and Anomalies and null values | We do not have any outliers and analomalies in our dataset except that we have null values so we have handle the null values and drop them in the way as follows |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data | ```
[7] df = pd.read_csv('measurements.csv')
    print(df.head())

      distance  consume  speed  temp_inside  temp_outside  specials gas_type  AC  \
   0      28.0      5.0     26         21.5            12       NaN      E10    0
   1      12.0      4.2     30         21.5            13       NaN      E10    0
   2      11.2      5.5     38         21.5            15       NaN      E10    0
   3      12.9      3.9     36         21.5            14       NaN      E10    0
   4      18.5      4.5     46         21.5            15       NaN      E10    0

      rain  sun  refill liters refill gas
   0     0    0           45.0         E10
   1     0    0            NaN         NaN
   2     0    0            NaN         NaN
   3     0    0            NaN         NaN
   4     0    0            NaN         NaN
``` |

| Handling Missing Data | df.isnull()<br><br>|  | distance | consume | speed | temp_inside | temp_outside | specials | gas_type | AC | rain | sun | refill liters | refill gas |<br>|---|---|---|---|---|---|---|---|---|---|---|---|---|<br>| 0 | False | False | False | False | False | True | False | False | False | False | False | False |<br>| 1 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 2 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 3 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 4 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |<br>| 383 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 384 | False | False | False | False | False | False | False | False | False | False | True | True |<br>| 385 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 386 | False | False | False | False | False | True | False | False | False | False | True | True |<br>| 387 | False | False | False | False | False | False | False | False | False | False | True | True |<br><br>388 rows × 12 columns |
| Data Transformation | **Splitting the data :**<br><br>```python<br>[24] x=x.values<br>     y=y.values<br>```<br><br>```python<br>#Splitting Data Into Train And Test<br>```<br>`+ Code`<br>```python<br>[26] x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=42)<br>``` |
| Feature Engineering | modifying existing ones<br><br>```python<br>[20] #seperating independent and dependent variables<br><br>[21] from sklearn.model_selection import train_test_split<br>     from sklearn.linear_model import LinearRegression<br><br>[22] x=df.drop(['consume','gas_type'],axis=1)<br>     y=df['consume']<br><br>[23] x.columns<br>     Index(['distance', 'speed', 'temp_inside', 'temp_outside', 'AC', 'rain',<br>            'sun'],<br>           dtype='object')<br><br>[24] x=x.values<br>     y=y.values<br>``` |
| Save Processed Data | Code to save the cleaned and processed data for future use.<br><br>```python<br>import pickle<br>pickle.dump(dt,open('fuel2.pkl','wb'))<br>``` |