

# Data Preprocessing and Visualization

## Vidita Dhavane

```
In [1]: #Import pandas for data manipulation
import pandas as pd

In [2]: #Load the dataset into a dataframe
dataFrame=pd.read_excel("Data Science - Intern - Data Set.xlsx")
dataFrame

Out [2]:
```

	uniqued	ts	lat	lng	external_bat_voltage	internal_bat_voltage	tanklevel	throttle	coolant	live_location	to_1	
0	11	220403501	1709340731	19.09340	78.355148	NaN	NaN	28	NaN	NaN	33 ...	NaN
1	11	220403501	1709340794	19.09347	78.355141	NaN	NaN	22	NaN	NaN	41 ...	NaN
2	11	220403501	1709340957	19.09344	78.355179	NaN	NaN	25	NaN	NaN	48 ...	NaN
3	11	220403501	1709340920	19.09328	78.355225	NaN	NaN	35	NaN	NaN	52 ...	NaN
4	11	220403501	1709340983	19.09394	78.355339	NaN	NaN	0	NaN	NaN	57 ...	NaN
...	...	...	...	...	...	...	...	...	...	...	...	...
18429	11	220403501	1713635965	18.869371	79.020393	NaN	NaN	50	NaN	NaN	89 ...	NaN
18430	11	220403501	1713639929	18.867130	79.018272	NaN	NaN	0	NaN	NaN	92 ...	NaN
18431	11	220403501	1713639993	18.862076	79.013924	NaN	NaN	10	NaN	NaN	87 ...	NaN
18432	11	220403501	1713639056	18.859493	79.011673	NaN	NaN	52	NaN	NaN	87 ...	NaN
18433	11	220403501	1713636120	18.858324	79.009033	NaN	NaN	72	NaN	NaN	90 ...	NaN

18434 rows x 141 columns

```
In [3]: #Overview of the dataset
dataFrame.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18434 entries, 0 to 18433
Columns: 141 entries, uniqued to adblue_level
dtypes: float64(122), int64(13), object(6)
memory usage: 19.8+ MB

In [4]: #See the datatypes for each column
dataFrame.dtypes

Out [4]:
```

	uniqued	ts	lat	lng	external_bat_voltage	boost_pressure	intake_manifold_1_temperature	air_inlet_pressure	exhaust_gas_temperature	adblue_level
	object	object	float64	float64	float64	float64	float64	float64	float64	float64
ts		int64								
lat		float64								
lng		float64								
external_bat_voltage		float64								
boost_pressure		float64								
intake_manifold_1_temperature		float64								
air_inlet_pressure		float64								
exhaust_gas_temperature		float64								
adblue_level		float64								
length	141, dtype: object									

```
In [5]: #Check the number of rows and columns
print("Number of rows in dataset are :",dataFrame.shape[0])
print("Number of columns in dataset are :",dataFrame.shape[1])

Number of rows in dataset are : 18434
Number of columns in dataset are : 141

In [6]: #check for null values
dataFrame.isnull().sum()

Out [6]:
```

	uniqued	ts	lat	lng	external_bat_voltage	boost_pressure	intake_manifold_1_temperature	air_inlet_pressure	exhaust_gas_temperature	adblue_level
	0	0	0	0	18434	18434	18434	18434	18434	3664
length	141, dtype: object									

## Data Cleaning

```
In [7]: #drop the entire column if it has null values or is empty
# Get the total number of rows
total_rows = len(dataFrame)

# Identify columns with all values as null
columns_to_drop = dataFrame.columns[dataFrame.isnull().sum() == total_rows]

# Drop those columns
dataFrame_cleaned = dataFrame.drop(columns=columns_to_drop)

print("Dropped columns:", columns_to_drop)

Dropped columns: Index(['external_bat_voltage', 'internal_bat_voltage', 'tanklevel', 'throttle',
                        'intakeairtemp', 'inasp', 'warmups', 'maf', 'fuelrail', 'mil',
                        ...,
                        'request_id', 'live_location', 'to_be_expired_at', 'status', 'event_id',
                        'particulatefiltertrap_inlet_pressure', 'boost_pressure',
                        'intake_manifold_1_temperature', 'air_inlet_pressure',
                        'exhaust_gas_temperature',
                        'exhaust_gas_temperature',
                        'exhaust_gas_temperature'], dtype=object)

In [8]: #check if the columns have been dropped successfully
dataFrame_cleaned

Out [8]:
```

	uniqued	ts	lat	lng	engineload	coolant	engineoiltemp	vehiclespeed	rpm	obddistance	...	engine_throttle_valve_pos	
0	11	220403501	1709340731	19.09340	78.355148	28	33	1774.9688	0.00000	749.000	60789000	—	100.0
1	11	220403501	1709340794	19.09347	78.355141	22	41	1774.9688	0.00000	1208.375	60789000	—	99.6
2	11	220403501	1709340957	19.09344	78.355179	25	48	1774.9688	1.40625	752.625	60789000	—	99.6
3	11	220403501	1709340920	19.09328	78.355225	35	52	1774.9688	2.81250	966.375	60789020	—	99.6
4	11	220403501	1709340983	19.09394	78.355339	0	57	1774.9688	4.53125	1065.750	60789040	—	100.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...
18429	11	220403501	1713635965	18.869371	79.020393	50	89	1774.9688	14.28125	945.125	71872976	—	100.0
18430	11	220403501	1713639929	18.867130	79.018272	0	92	1774.9688	38.26125	1130.750	71873344	—	100.0
18431	11	220403501	1713639993	18.862076	79.013924	10	87	1774.9688	32.34375	749.750	71874064	—	100.0
18432	11	220403501	1713639056	18.859493	79.011673	52	87	1774.9688	17.71875	820.625	71874464	—	100.0
18433	11	220403501	1713636120	18.858324	79.009033	72	90	1774.9688	16.21875	765.125	71874816	—	100.0

18434 rows x 33 columns

```
In [9]: #check if the previous step has been executed successfully
dataFrame_cleaned.isnull().sum()

Out [9]:
```

	uniqued	ts	lat	lng	engineload	coolant	engineoiltemp	vehiclespeed	rpm	obddistance	...	engine_throttle_valve_pos
	0	0	0	0	0	0	0	0	0	0	...	0
runtime											14	81
engine_torque_percent											0	0
selected_gear											0	0
current_gear											0	0
fuel_consumption											0	0
fuel_level											0	0
fuel_economy											0	0
accelerator_pedal_pos											0	0
pluscode											0	0
vibration_status											0	0
can_raw_data											0	0
engine_throttle_valve1_pos											0	0
engine_throttle_valve2_pos											0	0
drivers_demand_engine_torque_percent											0	0
engine_torque_mode											0	0
accelerator_pedal_pos_2											0	0
brake_switch_status											0	0
clutch_switch_status											0	0
parking_switch_status											0	0
adblue_level											3664	0
dtype:	int64											

```
In [10]: # Get the number of unique values in each column, if there is only 1 unique value in a column simply drop that as it w
on't provide any insights
unique_values_per_column = dataFrame_cleaned.nunique()
print(unique_values_per_column)

uniqued      1
ts           18434
lat          15255
lng           101
engineload   68
coolant       81
engineoiltemp 3692
vehiclespeed  3600
rpm           16747
obddistance  6195
engine_torque_percent 1
selected_gear 1
current_gear  20
fuel_consumption 6123
fuel_level     250
fuel_economy   251
accelerator_pedal_pos 251
pluscode      2439
vibration_status 1
can_raw_data  17874
engine_throttle_valve1_pos 112
engine_throttle_valve2_pos 1
drivers_demand_engine_torque_percent 83
engine_torque_mode 1
accelerator_pedal_pos_2 1
brake_switch_status 2
clutch_switch_status 2
parking_switch_status 2
adblue_level   131
dtype: int64

In [11]: # Count the number of zeros in each column
zero_counts = (dataFrame_cleaned == 0).sum()
print(zero_counts)

uniqued      0
ts           0
lat          0
lng          0
engineload   0
coolant       0
engineoiltemp 0
vehiclespeed   0
rpm           0
obddistance   0
engine_torque_percent 2210
selected_gear  0
current_gear  5340
fuel_consumption 0
fuel_level     0
fuel_economy   0
accelerator_pedal_pos 8066
pluscode      18434
vibration_status 1
can_raw_data  17874
engine_throttle_valve1_pos 112
engine_throttle_valve2_pos 1
drivers_demand_engine_torque_percent 83
engine_torque_mode 1
accelerator_pedal_pos_2 1
brake_switch_status 2
clutch_switch_status 2
parking_switch_status 2
adblue_level   0
dtype: int64
```

## Statistical Analysis

```
In [12]: #Import numpy as np
#replace zeros in lat and lng column with mean
columns_to_fill = ['lat', 'lng']
mean_values = dataFrame_cleaned[columns_to_fill].replace(0, np.nan).mean()
dataFrame_cleaned[columns_to_fill] = dataFrame_cleaned[columns_to_fill].replace(0, mean_values)

In [13]: # Calculate the median of non-zero values and fill column current_gear
median_gear = dataFrame_cleaned[dataFrame_cleaned['current_gear'] != 0]['current_gear'].median()
dataFrame_cleaned['current_gear'] = dataFrame_cleaned['current_gear'].replace(0, median_gear)

In [14]: # Calculate the mean of non-zero values
mean_engineload = dataFrame_cleaned[dataFrame_cleaned['engineload'] != 0]['engineload'].mean()
dataFrame_cleaned['engineload'] = dataFrame_cleaned['engineload'].replace(0, mean_engineload)

In [15]: # Find indices where 'fuel_economy' is 0
indices_to_drop = dataFrame_cleaned[dataFrame_cleaned['fuel_economy'] == 0].index
dataFrame_cleaned.drop(indices_to_drop, inplace=True)

In [16]: #drop all the column which have only 1 unique value
dataFrame_cleaned.drop(['uniqued','engineoiltemp','selected_gear','engine_throttle_valve2_pos','engi
nefueltemp','accelerator_pedal_pos_2','vibration_status'],axis='columns')

In [17]: print("Number of rows in dataset after cleaning are :", dataFrame_cleaned.shape[0])
print("Number of columns in dataset after cleaning are :", dataFrame_cleaned.shape[1])

Number of rows in dataset after cleaning are : 17354
Number of columns in dataset after cleaning are : 26

In [18]: #check for duplicates
duplicates = dataFrame_cleaned.duplicated()
print("Number of duplicate rows:", duplicates.sum())

Number of duplicate rows: 0

Overview of the statistical properties

In [19]: #check the overall statistics of the data
dataFrame_cleaned.describe()

Out [19]:
```

	ts	lat	lng	engineload	coolant	vehiclespeed	rpm	obddistance	runtime	engine_torque_percent	...	pluscode
count	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	1.735400e+04	...	1.735400e+04
mean	1.711640e+09	18.751160	78.992325	45.362505	85.032652	35.433739	948.176357	6.648797e+07	2132.320598	28.792440	...	28.792440
std	1.278986e+06	0.304881	0.386512	25.163278	7.043764	16.376169	333.256477	3.263905e+06	92.670827	22.866899	...	22.866899
min	1.709340e+09	18.717802	78.118062	1.000000	36.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000
50%	1.710904e+09	18.712358	78.761806	23.000000	84.000000	26.285000	767.125000	6.372202e+07	2053.162500	7.000000	...	7.000000
25%	1.711905e+09	18.848518	79.076946	45.021131	85.000000	39.500000	979.875000	6.652453e+07	2132.950000	25.500000	...	25.500000
75%	1.712165e+09	18.812683	79.132142	66.000000	89.000000	47.484375	1169.843750	6.940579e+07	2213.787500	49.000000	...	49.000000
max	1.713656e+09	19.096500	79.481766	100.000000	98.000000	81.703125	1975.250000	7.187482e+07	2287.100000	80.000000	...	80.000000

8 rows x 21 columns

## Feature Engineering

```
In [20]: #convert the ts column to date and time
from datetime import datetime

# Assuming you have a DataFrame named 'dataFrame_cleaned' with a column 'ts'
# Convert 'ts' from Unix timestamp to datetime
dataFrame_cleaned['ts_converted'] = pd.to_datetime(dataFrame_cleaned['ts'], unit='s')

# Print the first few rows to verify
print(dataFrame_cleaned[['ts', 'ts_converted']].head())

ts      ts_converted
2  1709340837  2024-03-02 00:54:17
3  1709340920  2024-03-02 00:55:20
4  1709340983  2024-03-02 00:56:13
5  1709341725  2024-03-02 01:08:45
6  1709344633  2024-03-02 01:57:13

In [21]: #extract the date only from date time format
dataFrame_cleaned['ts_converted'] = dataFrame_cleaned['ts_converted'].dt.date
dataFrame_cleaned['ts_converted'] = dataFrame_cleaned['ts_converted'].dt.strftime('%m/%d/%Y')

In [22]: dataFrame_cleaned

Out [22]:
```

	ts	lat	lng	engineload	coolant	vehiclespeed	rpm	obddistance	runtime	engine_torque_percent	...	pluscode
2	1709340837	19.090414	78.355179	25.000000	48	1.40625	752.625	60789000	1970.85	18	129	-7JFW9394+ OCF00X
3	1709340920	19.090258	78.355225	36.000000	52	2.81250	996.375	60789020	1970.90	27	129	-7JFW9394+ OCF00X
4	1709340983	19.090994	78.355339	45.021131	57	4.53125	1065.750	60789040	1970.90	1	3	-7JFW9394+ OCF00X
5	1709341725	19.091517	78.356079	22.000000	58	0.09375	750.750	60789020	1970.95	20	129	-7JFW9394+ OCF00X
6	1709344633	19.091517	78.356049	45.021131	53	0.00000	0.000	60789020	1970.95	0	16	-7JFW9394+ OCF00X
...	...	...	...	...	...	...	...	...	...	...	...	...
18429	1713635965	18.869371	79.020393	50.000000	89	14.28125	945.125	71872976	2287.00	33	11	-7JCKV29C+ OCF00C
18430	1713639929	18.867130	79.018272	45.021131	92	38.28125	1130.750	71873344	2287.05	0	16	-7JCKV29B+ OCF00C
18431	18.862076	18.862076	79.013924	10.000000	87	32.34375	749.750	71874064	2287.05	7	16	-7JCKV26T+ OCF00C
18432	18.859493	18.859493	79.011673	52.000000	80	17.71875	820.625	71874464	2287.05	39	13	-7JCKV25E+ OCF00C
18433	18.858324	18.858324	79.009033	72.000000	90	16.21875	765.125	71874816	2287.10	45	13	-7JCKV25S+ OCF00C

17354 rows x 27 columns

```
In [23]: # Identify categorical columns
categorical_columns = dataFrame_cleaned.select_dtypes(include=['object', 'category']).columns
print("Categorical columns:\n", categorical_columns)

Categorical columns:
Index(['pluscode', 'can_raw_data', 'brake_switch_status',
       'clutch_switch_status', 'parking_switch_status', 'ts_converted'],
      dtype='object')

In [24]: #check how many unique values are present in the specific column
unique_value_counts = dataFrame_cleaned['brake_switch_status'].value_counts()
print("Unique value counts :\n", unique_value_counts)

Unique value counts :
Released    10000
Pressed     1264
dtype: int64

Converting categorical values to numeric

In [25]: #converting categorical values to numeric where released=0 and pressed=1
# Define the mapping
mapping = {'Released': 0, 'Pressed': 1}

# List of columns to apply the mapping
columns_to_replace = ['brake_switch_status', 'clutch_switch_status', 'parking_switch_status']

# Replace values in the specified columns
dataFrame_cleaned[columns_to_replace] = dataFrame_cleaned[columns_to_replace].replace(mapping)

print("Updated columns:\n", dataFrame_cleaned[columns_to_replace].head())

Updated columns:
brake_switch_status clutch_switch_status parking_switch_status
0      Released      Released      Released
1      Released      Released      Released
2      Released      Released      Released
3      Released      Released      Released
4      Released      Released      Released

In [26]: dataFrame_cleaned=dataFrame_cleaned.drop(columns=['ts'])
dataFrame_cleaned

Out [26]:
```

	lat	lng	engineload	coolant	vehiclespeed	rpm	obddistance	runtime	engine_torque_percent	current_gear	...	pluscode
2	19.090414	78.355179	25.000000	48	1.40625	752.625	60789000	1970.85	18	129	-7JFW9394+ OCF00X	
3	19.090258	78.355225	36.000000	52	2.81250	996.375	60789020	1970.90	27	129	-7JFW9394+ OCF00X	
4	19.090994	78.355339	45.021131	57	4.53125	1065.750	60789040	1970.90	1	3	-7JFW9394+ OCF00X	
5	19.091517	78.356079	22.000000	58	0.09375	750.750	60789020	1970.95	20	129	-7JFW9394+ OCF00X	
6	19.091517	78.356										



# Report

## 1. Data loading and preparation:

- Data loading and preparation done which included identification of variables , determining the size of dataset ,identifying missing values etc

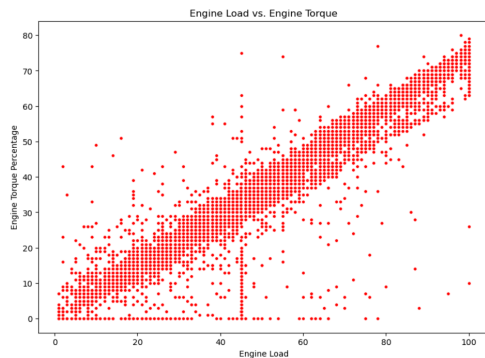
## 2. Data cleaning and preprocessing:

- Handled missing values
- Converted data type . Converted the columns brake\_switch\_status , clutch\_switch\_status and parking\_switch\_status values to numeric values 0 and 1 . The values were initially categorical i.e pressed and released .
- Handled duplicate values

## 3. Exploratory data analysis:

- Prepared visuals for data analysis like scatter plots, column charts etc

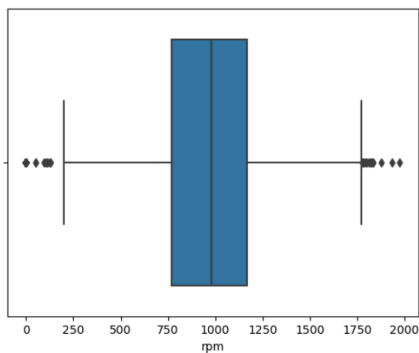
### 1. Analyzing vehicle performance and efficiency by using the below two columns Engine Load and Torque



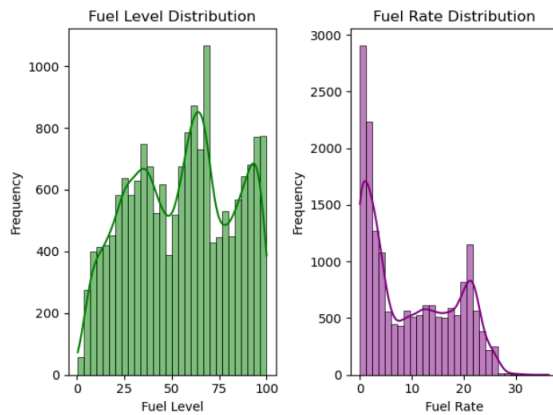
- ### 2. The scatter plot visualizes the relationship between vehicle speed and RPM. It helps identify any correlation or pattern between these two variables, such as whether higher speeds correspond to higher RPMs or if there's a consistent trend.



### 3. Detecting outliers in Engine Data using boxplots which could indicate mechanical issues or unusual driving behavior. i.e outliers

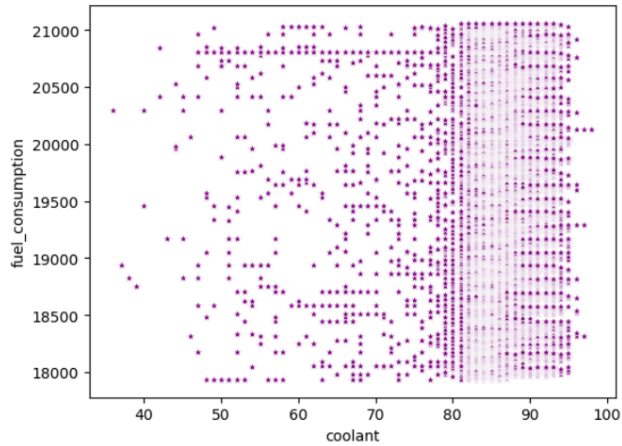


### 4. Distribution of fuel\_level and fuel\_rate

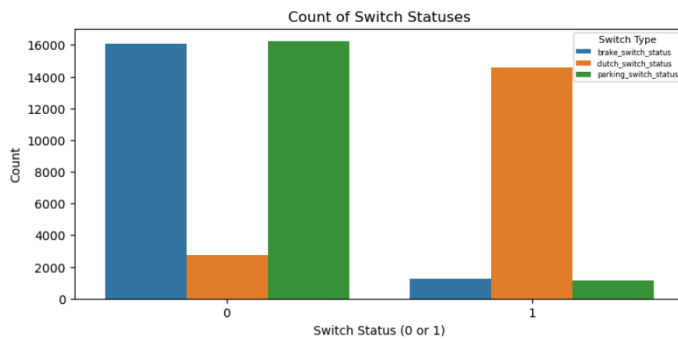


### 5. Impact of Environmental Factors

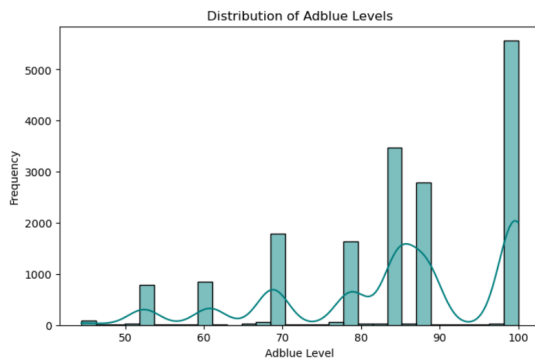
**Coolant Temperature and Engine Performance: Exploring the relationship between coolant temperature and fuel\_consumption**



6. The bar chart shows the frequency of 'Pressed' (1) and 'Released' (0) statuses for the brake, clutch, and parking switches, revealing their activity patterns. It helps compare the usage or engagement levels of each switch type.



7. The distribution, trends, and relationships of Adblue Level in the dataset.



#### 4. Feature Engineering:

- Created a new feature `ts_converted` where the timestamp column `ts` was converted to a proper date and time format .

## 5. Insights and Recommendations:

### Insights from the data:

#### 1. Fuel Efficiency:

- **Fuel Consumption and Coolant:** A higher coolant temperature might indicate more engine strain, leading to increased fuel consumption.
- **Optimal RPM and Speed:** Keeping the vehicle within a moderate RPM and speed range improves fuel economy.

#### 2. Vehicle Performance:

- **RPM and Torque Correlation:** High RPM results in increased torque, but excessive RPM can reduce fuel efficiency.
- **Current Gear and Speed:** Higher gears at appropriate speeds optimize both fuel consumption and vehicle performance.

#### 3. Safety:

- **Brake Usage:** Frequent brake usage may indicate aggressive driving, which can reduce fuel efficiency and lead to quicker wear on brake components.

#### 4. Emissions:

- **Adblue Level:** Maintaining optimal levels can help in reducing NOx emissions, crucial for environmental regulations compliance.

### Recommendations:

#### 1. Improving Fuel Efficiency:

- Encourage drivers to maintain a moderate RPM (e.g., between 2,000 and 3,000) for optimal fuel consumption.
- Educate drivers about coasting and accelerating smoothly to minimize sudden brake usage, which can save fuel.

#### 2. Enhancing Vehicle Performance:

- Encourage timely gear shifting and avoid running the engine at too high an RPM for prolonged periods.
- Regularly monitor coolant temperature and engine load to prevent engine overheating, which can strain performance.

#### 3. Boosting Safety:

- Introduce driver training on defensive driving techniques, such as minimizing aggressive braking and acceleration.
- Implement automated warnings when clutch or brake switches are pressed excessively.

#### 4. Reducing Emissions:

- Regularly check and maintain AdBlue levels in vehicles to ensure compliance with emission standards.
- Monitor and replace worn components that affect engine load and fuel consumption, reducing emissions further.

Below is the PowerBI dashboard for better visualization.

