# Python Extract Transform and Load Project

## Introduction

In this project, I analyse fuel economy between different car manufacturers for an Automobile research company, that wants to launch new research based on the given data.

## ETL process

Extract:

- Downloaded raw data files (CSV or XLSX) from the data source ( 📊 Download Fuel Economy Data )
- I used Python to read these files and extract the data into a Pandas data frame.

Transform:

- I performed data cleaning operations, including missing values, renaming columns, and dropping duplicates.
- I calculated additional metrics such as checksums for data integrity.
- I visualized trends, distributions, and relationships in the data, which can be considered a form of transformation for analysis purposes.

Load:

- I cleaned and transformed data into a new CSV (called 'cleaned_data.csv') for the main analysis.
- I aggregated data to answer specific questions.
- I stored the results of all ten visualisations for analysis.

## Python - Data cleaning process

∨ Data cleaning.py

```python
import os
import pandas as pd
import glob

def clean_and_concatenate_data(folder_path):
    # Step 1: Read Excel files using the specified folder path
    excel_files = glob.glob(os.path.join(folder_path, '*.xlsx'))

    # Initialize an empty list to store DataFrames
    dfs = []

    # Read each Excel file and append to the list of DataFrames
    for file in excel_files:
        try:
            df = pd.read_excel(file)

            # Append the DataFrame to the list
```
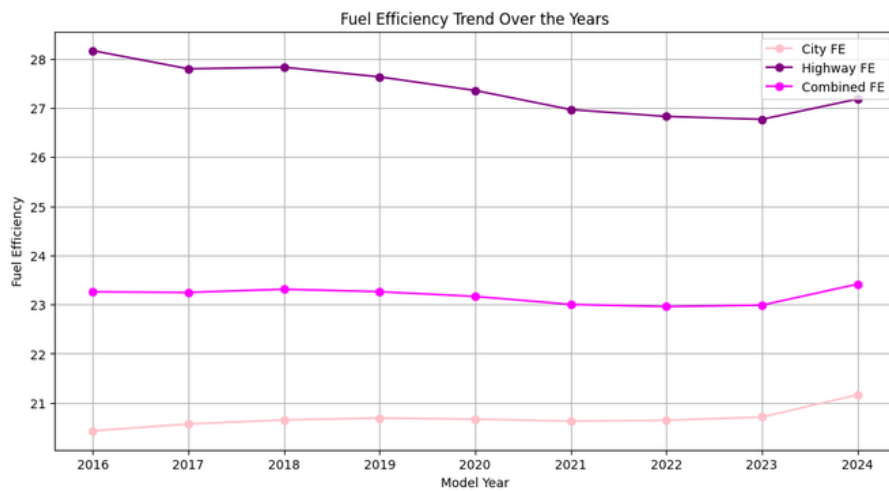
```python
18            dfs.append(df)
19
20        except pd.errors.ParserError as e:
21            print(f"Error reading file {file}: {e}")
22
23    # Step 2: Concatenate all DataFrames into a single DataFrame
24    concatenated_df = pd.concat(dfs, ignore_index=True)
25
26    # Step 3: Clean the DataFrame
27    # Rename columns, drop duplicates, handle missing values, etc.
28    # Define columns to be renamed
29    columns_to_rename = {
30        'Eng Displ': 'Engine Displacement',
31        '# Cyl': '# Cylinders',
32        'City FE (Guide) - Conventional Fuel': 'City FE',
33        'Hwy FE (Guide) - Conventional Fuel': 'Highway FE',
34        'Comb FE (Guide) - Conventional Fuel': 'Combined FE',
35        'Air Aspiration Method Desc': 'Air Aspiration Method',
36        'Trans Desc': 'Transmission Description',
37        'City CO2 Rounded Adjusted': 'City CO2',
38        'Hwy CO2 Rounded Adjusted': 'Highway CO2',
39        'Comb CO2 Rounded Adjusted (as shown on FE Label)': 'Combined CO2'
40    }
41
42    # Columns to retain
43    columns_to_retain = [
44        'Model Year', 'Mfr Name', 'Division', 'Carline', 'Engine Displacement', '# Cylinders', 'Transmission
45        'City FE', 'Highway FE', 'Combined FE', 'Air Aspiration Method', 'Transmission Description', '# Gears
46        'Drive Desc', 'Carline Class Desc', 'Release Date', 'City CO2', 'Highway CO2', 'Combined CO2'
47    ]
48
49    # Rename columns
50    concatenated_df.rename(columns=columns_to_rename, inplace=True)
51
52    # Retain specified columns and drop any other columns
53    concatenated_df = concatenated_df[columns_to_retain]
54
55    # Replace null values with None
56    concatenated_df = concatenated_df.where(pd.notna(concatenated_df), None)
57
58    # Remove duplicate columns
59    concatenated_df = concatenated_df.loc[:, ~concatenated_df.columns.duplicated()]
60
61    # Step 4: Define the output file path
62    output_file_path = os.path.join(folder_path, 'cleaned_data.csv')
63
64    # Step 5: Save the cleaned DataFrame to a new CSV file
65    concatenated_df.to_csv(output_file_path, index=False, encoding='utf-8')
66
67    # Step 6: Print a message indicating successful completion
68    print(f"Data cleaning and concatenation completed. Cleaned data saved to: {output_file_path}")
69
70    # Step 7: Return the cleaned DataFrame
71    return concatenated_df
72
73 # Call the function with the specified folder path and store the result
74 folder_path = folder_path = "C:\\Users\\Esther Oluwaseye\\Desktop\\Git_assessment_local\\python-data-assessme
75 cleaned_data = clean_and_concatenate_data(folder_path)
```
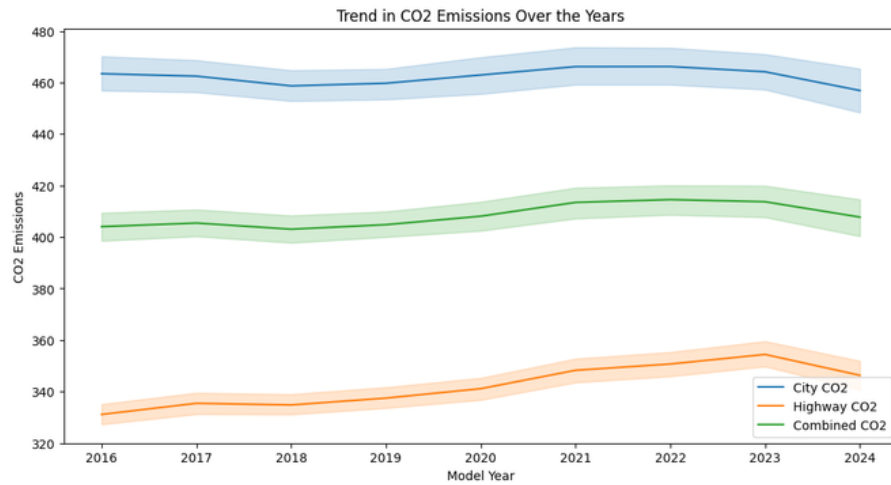
**Python data analysis questions**

```python
1   #Question 1: How has fuel efficiency changed over the years
2
3   import pandas as pd
4   import matplotlib.pyplot as plt
5
6   # Load cleaned data frame
7   cleaned_data = pd.read_csv("cleaned_data.csv")
8
9   # Group the data by 'Model Year' and calculate the mean for each year
10  fuel_efficiency_trend = cleaned_data.groupby('Model Year')[['City FE', 'Highway FE', 'Combined FE']].mean()
11
12  # Plotting the trend with pink and purple colors
13  plt.figure(figsize=(12, 6))
14  plt.plot(fuel_efficiency_trend.index, fuel_efficiency_trend['City FE'], label='City FE', marker='o', color='pink
15  plt.plot(fuel_efficiency_trend.index, fuel_efficiency_trend['Highway FE'], label='Highway FE', marker='o', color
16  plt.plot(fuel_efficiency_trend.index, fuel_efficiency_trend['Combined FE'], label='Combined FE', marker='o', col
17
18  # Adding labels and title
19  plt.xlabel('Model Year')
20  plt.ylabel('Fuel Efficiency')
21  plt.title('Fuel Efficiency Trend Over the Years')
22  plt.legend()
23  plt.grid(True)
24  plt.show()
25
```

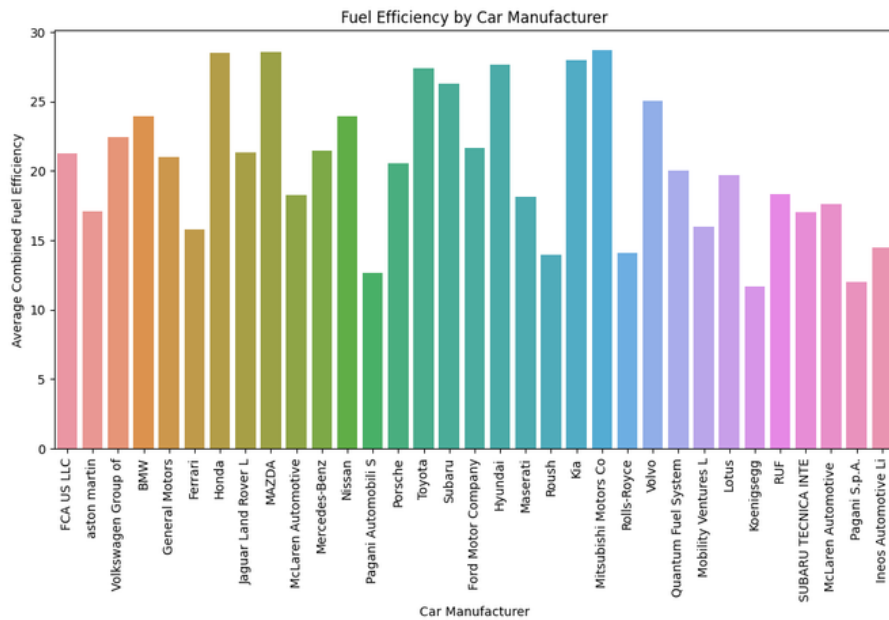Question 2: Look at the trend of CO2 emissions over the years

Trend in CO2 Emissions Over the Years

```
1   # Q2:Look at the trend of CO2 emissions over the years
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   # Set the size of the plot
7   plt.figure(figsize=(12, 6))
8
9   # Plot the trend in CO2 emissions over the years for City, Highway, and Combined
10  sns.lineplot(x='Model Year', y='City CO2', data=cleaned_data, label='City CO2')
11  sns.lineplot(x='Model Year', y='Highway CO2', data=cleaned_data, label='Highway CO2')
12  sns.lineplot(x='Model Year', y='Combined CO2', data=cleaned_data, label='Combined CO2')
13
14  # Label the x-axis
15  plt.xlabel('Model Year')
16
17  # Label the y-axis
18  plt.ylabel('CO2 Emissions')
19
20  # Add a title to the plot
21  plt.title('Trend in CO2 Emissions Over the Years')
22
23  # Display legend to distinguish between City, Highway, and Combined CO2
24  plt.legend()
25
26  # Show the plot
27  plt.show()
28
```

Q3: Which car manufacturers produce the most fuel-efficient vehicles?

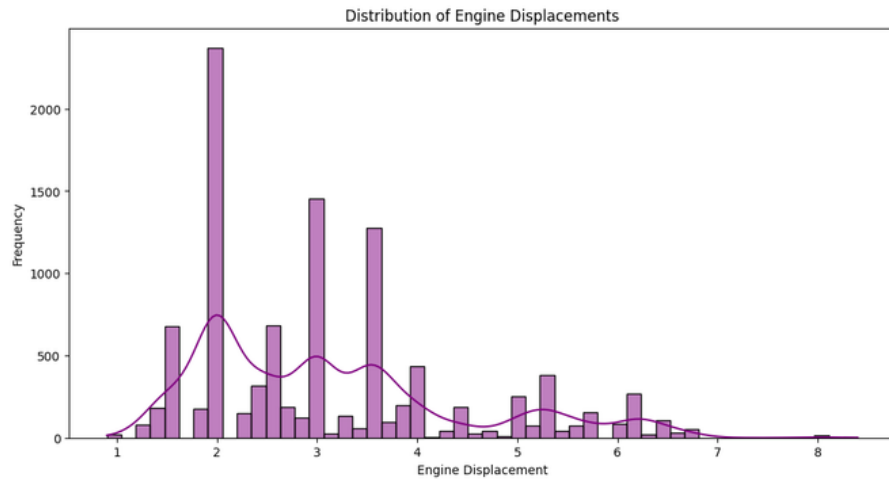Fuel Efficiency by Car Manufacturer

```
1   # Q3: Which car manufacturers produce the most fuel-efficient vehicles?
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   # Set the size of the plot
7   plt.figure(figsize=(12, 6))
8
9   # Create a bar plot for average combined fuel efficiency by car manufacturer
10  sns.barplot(x='Mfr Name', y='Combined FE', data=cleaned_data, errorbar=None, estimator='mean')
11
12  # Rotate x-axis labels for better visibility
13  plt.xticks(rotation=90)
14
15  # Label the x-axis
16  plt.xlabel('Car Manufacturer')
17
18  # Label the y-axis
19  plt.ylabel('Average Combined Fuel Efficiency')
20
21  # Add a title to the plot
22  plt.title('Fuel Efficiency by Car Manufacturer')
23
24  # Show the plot
25  plt.show()
26
```

Q4: What is the distribution of engine displacements among different car models

Distribution of Engine Displacements

```
1   # Q4:What is the distribution of engine displacements among different car models?
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   # Set the size of the plot
7   plt.figure(figsize=(12, 6))
8
9   # Create a histogram plot with KDE for the distribution of engine displacements
10  sns.histplot(cleaned_data['Engine Displacement'], kde=True, color='purple')
11
12  # Label the x-axis
13  plt.xlabel('Engine Displacement')
14
15  # Label the y-axis
16  plt.ylabel('Frequency')
17
18  # Add a title to the plot
19  plt.title('Distribution of Engine Displacements')
20
21  # Show the plot
22  plt.show()
23
24
```

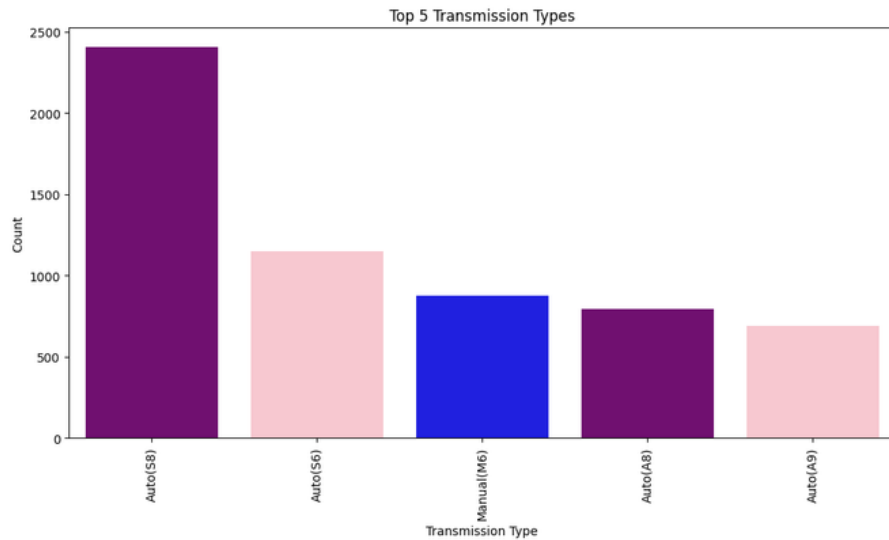Q5: How do different car classes compare in terms of fuel efficiency?

## Fuel Efficiency by Car Class



```python
1   # Q5:How do different car classes compare in terms of fuel efficiency?
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   # Set the size of the plot
7   plt.figure(figsize=(12, 6))
8
9   # Create a bar plot for the average combined fuel efficiency across different car classes
10  sns.barplot(x='Carline Class Desc', y='Combined FE', data=cleaned_data, errorbar=None, estimator='mean')
11
12  # Rotate x-axis labels for better readability
13  plt.xticks(rotation=90)
14
15  # Label the x-axis
16  plt.xlabel('Car Class')
17
18  # Label the y-axis
19  plt.ylabel('Average Combined Fuel Efficiency')
20
21  # Add a title to the plot
22  plt.title('Fuel Efficiency by Car Class')
23
24  # Show the plot
25  plt.show()
26
```

Q6: Are there any notable trends in the transmission types used in vehicles?
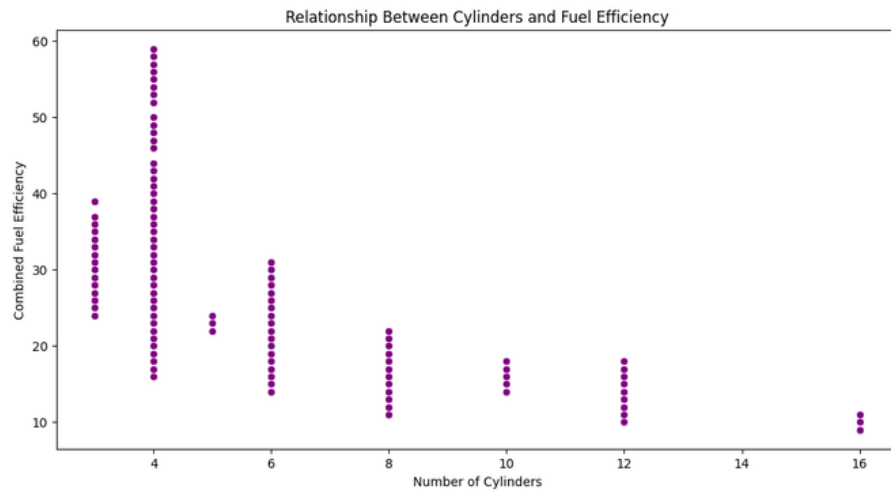
Top 5 Transmission Types

```python
1   # Q6: Are there any notable trends in the transmission types used in vehicles?
2   import matplotlib.pyplot as plt
3   import seaborn as sns
4
5   # Get the top 5 transmission types
6   top_transmissions = cleaned_data['Transmission'].value_counts().nlargest(5).index
7
8   # Filter the DataFrame for the top 5 transmission types
9   filtered_data = cleaned_data[cleaned_data['Transmission'].isin(top_transmissions)]
10
11  # Set the size of the plot
12  plt.figure(figsize=(12, 6))
13
14  # Create a count plot to visualize the distribution of transmission types
15  sns.countplot(x='Transmission', data=filtered_data, order=top_transmissions, palette=['purple', 'pink','blue'])
16
17  # Rotate x-axis labels for better readability
18  plt.xticks(rotation=90)
19
20  # Label the x-axis
21  plt.xlabel('Transmission Type')
22
23  # Label the y-axis
24  plt.ylabel('Count')
25
26  # Add a title to the plot
27  plt.title('Top 5 Transmission Types')
28
29  # Show the plot
30  plt.show()
31
32
```

Q7: What is the relationship between the number of cylinders and fuel efficiency?

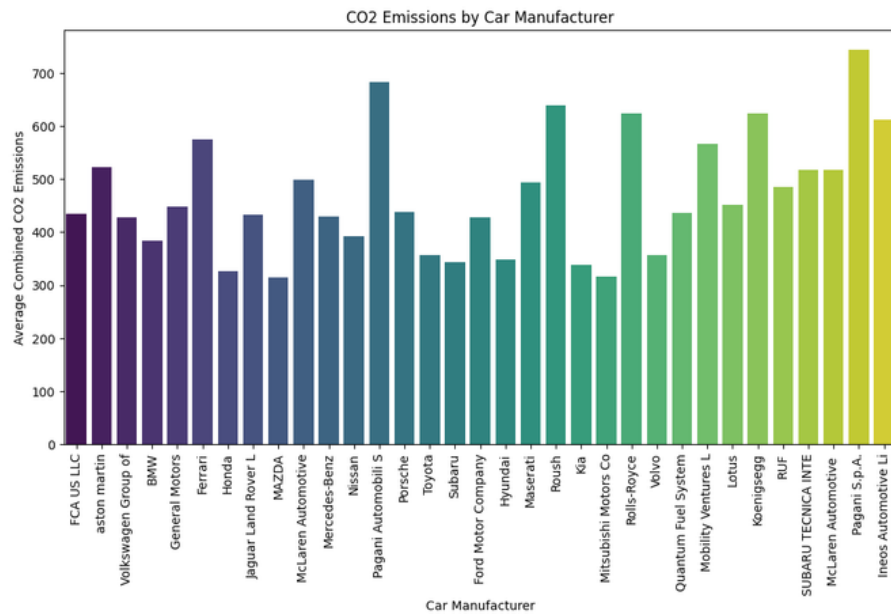Relationship Between Cylinders and Fuel Efficiency

```
1   import matplotlib.pyplot as plt
2   import seaborn as sns
3
4   # Set the size of the plot
5   plt.figure(figsize=(12, 6))
6
7   # Create a scatter plot to visualize the relationship
8   sns.scatterplot(x='# Cylinders', y='Combined FE', data=cleaned_data, color='purple')
9
10  # Label the x-axis
11  plt.xlabel('Number of Cylinders')
12
13  # Label the y-axis
14  plt.ylabel('Combined Fuel Efficiency')
15
16  # Add a title to the plot
17  plt.title('Relationship Between Cylinders and Fuel Efficiency')
18
19  # Show the plot
20  plt.show()
21
```

Q8: Which car manufacturers produce the lowest CO2 emissions?

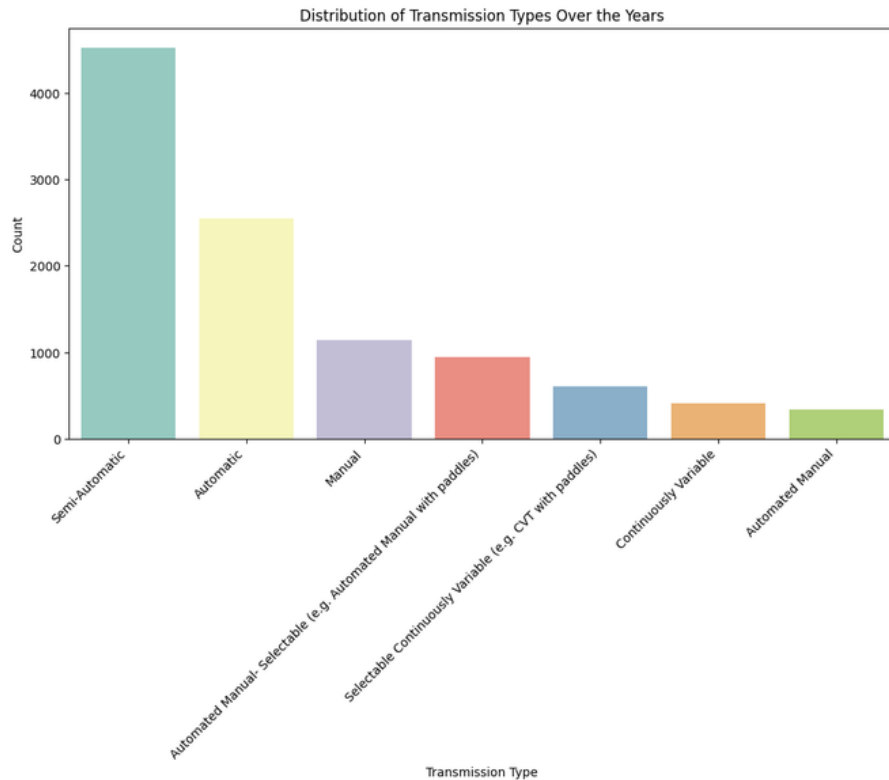CO2 Emissions by Car Manufacturer

```
1   # Q8: Which car manufacturers produce the lowest CO2 emissions?
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   # Set the size of the plot
7   plt.figure(figsize=(12, 6))
8
9   # Create a bar plot for the average combined CO2 emissions by car manufacturer
10  sns.barplot(x='Mfr Name', y='Combined CO2', data=cleaned_data, errorbar=None, estimator='mean', palette='viridis
11
12  # Rotate x-axis labels for better readability
13  plt.xticks(rotation=90)
14
15  # Label the x-axis
16  plt.xlabel('Car Manufacturer')
17
18  # Label the y-axis
19  plt.ylabel('Average Combined CO2 Emissions')
20
21  # Add a title to the plot
22  plt.title('CO2 Emissions by Car Manufacturer')
23
24  # Show the plot
25  plt.show()
26
27
```
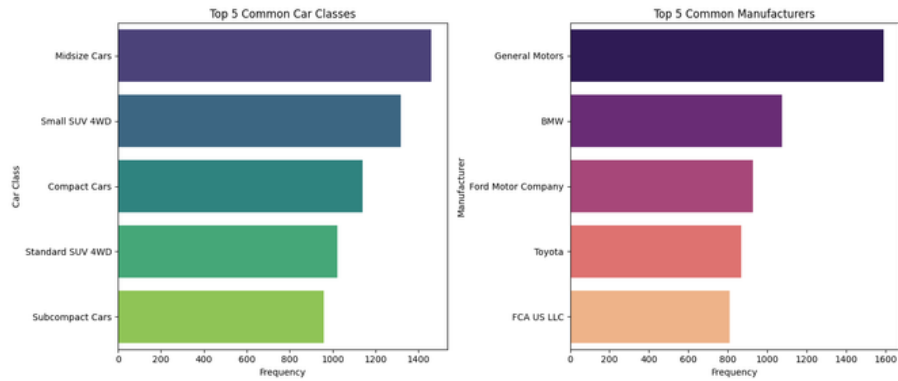
Q9: What are the top distribution of transmission types

Distribution of Transmission Types Over the Years



```
1   import matplotlib.pyplot as plt
2   import seaborn as sns
3
4   # Set the size of the plot
5   plt.figure(figsize=(12, 6))
6
7   # Create a count plot for the distribution of transmission types over the years
8   sns.countplot(x='Transmission Description', data=cleaned_data, palette='Set3', order=cleaned_data['Transmission
9
10  # Add a title to the plot
11  plt.title('Distribution of Transmission Types Over the Years')
12
13  # Label the x-axis
14  plt.xlabel('Transmission Type')
15
16  # Label the y-axis
17  plt.ylabel('Count')
18
19  # Rotate x-axis labels for better readability
20  plt.xticks(rotation=45, ha='right')
21
22  # Show the plot
23  plt.show()
24
```

Q10: What are the top 5 common car classes and manufacturers?

Top 5 Common Car Classes — Top 5 Common Manufacturers

```python
# Q10: What are the top 5 common car classes and manufacturers?

import matplotlib.pyplot as plt
import seaborn as sns

# Top 5 common car classes
top_car_classes = cleaned_data['Carline Class Desc'].value_counts().head(5)

# Top 5 common manufacturers
top_manufacturers = cleaned_data['Mfr Name'].value_counts().head(5)

# Plotting the bar plots
plt.figure(figsize=(14, 6))

# Plot for car classes
plt.subplot(1, 2, 1)
sns.barplot(x=top_car_classes.values, y=top_car_classes.index, palette='viridis')
plt.title('Top 5 Common Car Classes')
plt.xlabel('Frequency')
plt.ylabel('Car Class')

# Plot for manufacturers
plt.subplot(1, 2, 2)
sns.barplot(x=top_manufacturers.values, y=top_manufacturers.index, palette='magma')
plt.title('Top 5 Common Manufacturers')
plt.xlabel('Frequency')
plt.ylabel('Manufacturer')

plt.tight_layout()
plt.show()
```