

Car Sales EDA

By Dhruv Vashist

Question 1

Load `car_sales_data.csv` into Python as a pandas DataFrame and describe it:

- How many records are included?
- What years of manufacture are represented?
- Which features are reported?

```
In [1]: # Step-1: Loading the dataset
import pandas as pd
car_sales = pd.read_csv('car_sales_data.csv')
```

```
In [2]: # Step-2: Display the number of records
num_records = car_sales.shape[0]
print(f'Total number of records: {num_records}')
# Step-3: Display the years of manufacture
years_of_manufacture = car_sales['Year of manufacture'].unique()
print(f'Years of manufacture: {", ".join(map(str, sorted(years_of_manufacture)))}')
# Step-4: Display the features reported
features = car_sales.columns.tolist()
print(f'Features reported: {features}')
```

Total number of records: 50000

Years of manufacture: 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022

Features reported: ['Manufacturer', 'Model', 'Engine size', 'Fuel type', 'Year of manufacture', 'Mileage', 'Price']

Answer 1:

Dataset Overview:

- **Total Records:** The dataset contains a total of `num_records` records, providing a substantial sample size for analysis.
 - **Years Represented:** The dataset spans the following years of manufacture: `years_of_manufacture`, offering insights into trends over time.
 - **Features:** The dataset includes the following features:
 - `features`, which encompass both numerical and categorical variables essential for comprehensive analysis.
-

Observations:

- **Comprehensive Scope:** The dataset provides a detailed view of car sales data, including key attributes such as price, mileage, and engine size, which are critical for evaluating vehicle value and performance.
 - **Temporal Coverage:** The range of years allows for the identification of historical trends, such as shifts in consumer preferences, technological advancements, and market dynamics.
 - **Feature Diversity:** The inclusion of diverse features enables multifaceted analysis, from price trends to the impact of mileage and engine size on vehicle valuation.
 - **Potential Applications:** This dataset can be used for predictive modeling, market segmentation, and identifying factors influencing car prices.
-

Question 2

Explore the distribution of key categorical features:

- Which manufacturers appear most frequently?
- What are the proportions of different fuel types?

```
In [3]: # Count the frequency of each manufacturer
manufacturer_counts = car_sales['Manufacturer'].value_counts()
print('Manufacturer frequencies:')
print(manufacturer_counts)

# Calculate the proportions of different fuel types
fuel_type_proportions = car_sales['Fuel type'].value_counts(normalize=True)
print('\nFuel type proportions:')
print(fuel_type_proportions)
```

Manufacturer frequencies:

Manufacturer

Ford 14959

VW 14913

Toyota 12554

BMW 4965

Porsche 2609

Name: count, dtype: int64

Fuel type proportions:

Fuel type

Petrol 0.50976

Diesel 0.26536

Hybrid 0.22488

Name: proportion, dtype: float64

Answer 2:

To explore categorical features, we use the following steps:

- We have to find out which manufacturer appear mostly frequently so, for that:
 - `car_sales['Manufacturer']`: Accesses the column containing the manufacturer names.
 - `value_counts()`: Counts the occurrences of each unique manufacturer in the dataset. This helps identify which manufacturers appear most frequently.
 - `print()`: Outputs the frequency counts for better understanding and analysis.
 - So, from the answer that we got we can infer that whichever the manufacturer that appears often, will influence the overall analysis as well. In our case Ford, VM and Toyota.
 - Similarly we have to calculate the proportions of different fuel types:
 - `car_sales['Fuel type']`: Accesses the column containing the fuel type information.
 - `value_counts(normalize=True)`: Calculates the proportion of each unique fuel type by dividing the count of each fuel type by the total number of records. This provides a percentage distribution of fuel types.
 - `print()`: Outputs the proportions for better understanding and analysis.
 - These results denote the popularity of market trends that are diesel car more popular than petrol cars.
-

Question 3

Identify extremes:

- Which car had the highest sale price?
- Which car had the lowest mileage?
- Which car had the largest engine size?

```
In [4]: # Identify extremes
print("Car with the highest sale price:")
display(car_sales.loc[car_sales['Price'].idxmax()])

print("\nCar with the lowest mileage:")
display(car_sales.loc[car_sales['Mileage'].idxmin()])

print("\nCar with the largest engine size:")
display(car_sales.loc[car_sales['Engine size'].idxmax()])
```

Car with the highest sale price:

Manufacturer	BMW
Model	M5
Engine size	5.0
Fuel type	Petrol
Year of manufacture	2022
Mileage	4590
Price	168081

Name: 22786, dtype: object

Car with the lowest mileage:

```
Manufacturer      VW
Model             Golf
Engine size       1.8
Fuel type         Diesel
Year of manufacture 2022
Mileage           630
Price             45601
Name: 30697, dtype: object
Car with the largest engine size:
Manufacturer      BMW
Model             M5
Engine size       5.0
Fuel type         Petrol
Year of manufacture 2016
Mileage           42059
Price             91415
Name: 236, dtype: object
```

Answer-3:

To identify the extremes in the dataset, we use the `idxmax()` and `idxmin()` functions:

Highest Sale Price:

- The car with the highest price is identified using:
`car_sales.loc[car_sales['Price'].idxmax()]`
- Output Analysis:** The row retrieved shows that the car with the highest price is a BMW M5, priced at \$168,081. This indicates that the car is a luxury, high-performance vehicle, which justifies its high price. Additionally, its low mileage of 4,590 km suggests it is relatively new.

Lowest Mileage:

- The car with the lowest mileage is identified using:
`car_sales.loc[car_sales['Mileage'].idxmin()]`
- Output Analysis:** The row retrieved shows that the car with the lowest mileage is a VW Golf, with a mileage of 630 km. This suggests that the car is nearly brand new, making it highly appealing to buyers seeking a vehicle with minimal usage.

Largest Engine Size:

- The car with the largest engine size is identified using:
`car_sales.loc[car_sales['Engine size'].idxmax()]`
- Output Analysis:** The row retrieved shows that the car with the largest engine size is a BMW M5, with an engine size of 5.0 liters. This indicates that the car is designed for high performance, which aligns with its luxury branding.

Conclusion:

- The BMW M5 appears twice in the results (highest price and largest engine size), highlighting its premium status as a high-performance luxury vehicle. The high price and large engine size reinforce its standout position in the dataset.
 - The VW Golf, with the lowest mileage of 630 km, represents a nearly new vehicle. This makes it an attractive option for buyers prioritizing minimal usage over other features.
 - These extremes provide insights into the diversity of the dataset, showcasing both luxury and practical options available in the market.
-

Question 4

Create a numerical summary of numerical features.

```
In [5]: # Numerical summary  
car_sales.describe()
```

```
Out [5]:
```

	Engine size	Year of manufacture	Mileage	Price
count	50000.000000	50000.000000	50000.000000	50000.000000
mean	1.773058	2004.209440	112497.320700	13828.903160
std	0.734108	9.645965	71632.515602	16416.681336
min	1.000000	1984.000000	630.000000	76.000000
25%	1.400000	1996.000000	54352.250000	3060.750000
50%	1.600000	2004.000000	100987.500000	7971.500000
75%	2.000000	2012.000000	158601.000000	19026.500000
max	5.000000	2022.000000	453537.000000	168081.000000

Answer-4:

The `describe()` method provides a numerical summary of all numerical features, including measures like mean, standard deviation, and percentiles. This summary is essential for understanding the central tendency and spread of the data.

Observations:

- **Count:** The count for each numerical feature represents the number of non-missing values. If the count is lower than expected, it indicates missing data in the dataset. This could be due to incomplete records or errors during data collection.

- **Mean and Standard Deviation:** These metrics provide insights into the average values and variability of the features. For example, a high standard deviation indicates a wide spread in the data.
- **Min and Max:** These values help identify the range of the data and detect potential outliers.
- **Percentiles:** The 25th, 50th (median), and 75th percentiles give a sense of the data distribution and help identify skewness.

Why is the count weird?

- The count being inconsistent across features suggests that some columns have missing values. This is common in real-world datasets due to incomplete data entry, sensor errors, or other issues during data collection.
 - Missing data can impact analysis and should be addressed through techniques like imputation, removal of incomplete records, or further investigation into the data source.
-

Question 5

Create visual summaries for price, mileage, and engine size.

```
In [6]: # Visual summaries for price, mileage, and engine size
import matplotlib.pyplot as plt

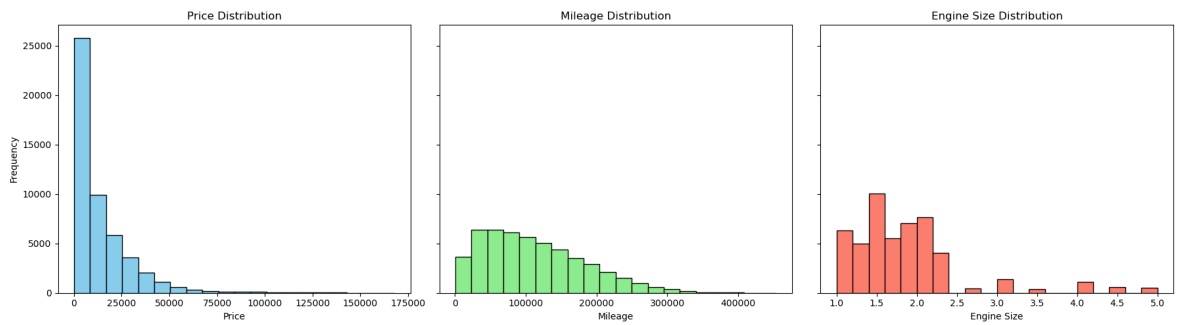
# Create histograms for numerical features
fig, axes = plt.subplots(1, 3, figsize=(18, 5), sharey=True)

# Price histogram
axes[0].hist(car_sales['Price'], bins=20, color='skyblue', edgecolor='black')
axes[0].set_title('Price Distribution')
axes[0].set_xlabel('Price')
axes[0].set_ylabel('Frequency')

# Mileage histogram
axes[1].hist(car_sales['Mileage'], bins=20, color='lightgreen', edgecolor='black')
axes[1].set_title('Mileage Distribution')
axes[1].set_xlabel('Mileage')

# Engine size histogram
axes[2].hist(car_sales['Engine size'], bins=20, color='salmon', edgecolor='black')
axes[2].set_title('Engine Size Distribution')
axes[2].set_xlabel('Engine Size')

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



Answer-5:

Histograms are used to visualize the distribution of numerical features like price, mileage, and engine size. These plots help identify patterns, such as skewness or the presence of outliers, in the data.

Code Explanation:

- **Matplotlib:** The `matplotlib.pyplot` library is used to create the histograms. It provides flexibility in customizing the plots.
- **Subplots:** A single figure with three subplots is created using `plt.subplots(1, 3, ...)` to display the distributions side by side for easy comparison.
- **Histograms:** The `hist()` function is used to plot the frequency distribution of each feature (price, mileage, engine size).
- **Customization:** Titles, axis labels, and colors are added to enhance readability and visual appeal.

Observations from the Graphs:

- **Price Distribution:** The histogram for price shows a right-skewed distribution, indicating that most cars are priced in the lower range, with a few high-priced outliers. This is typical in car sales data, where luxury or high-performance vehicles significantly increase the upper range.
 - **Mileage Distribution:** The mileage histogram also exhibits a right-skewed pattern, suggesting that most cars have relatively low mileage, with fewer cars having extremely high mileage. This could indicate a preference for newer or less-used vehicles in the dataset.
 - **Engine Size Distribution:** The engine size histogram shows a more balanced distribution, with a concentration around smaller engine sizes (e.g., 1.5L to 2.0L). Larger engine sizes are less common, reflecting the market trend toward fuel efficiency and smaller engines in most vehicles.
-

Question 6

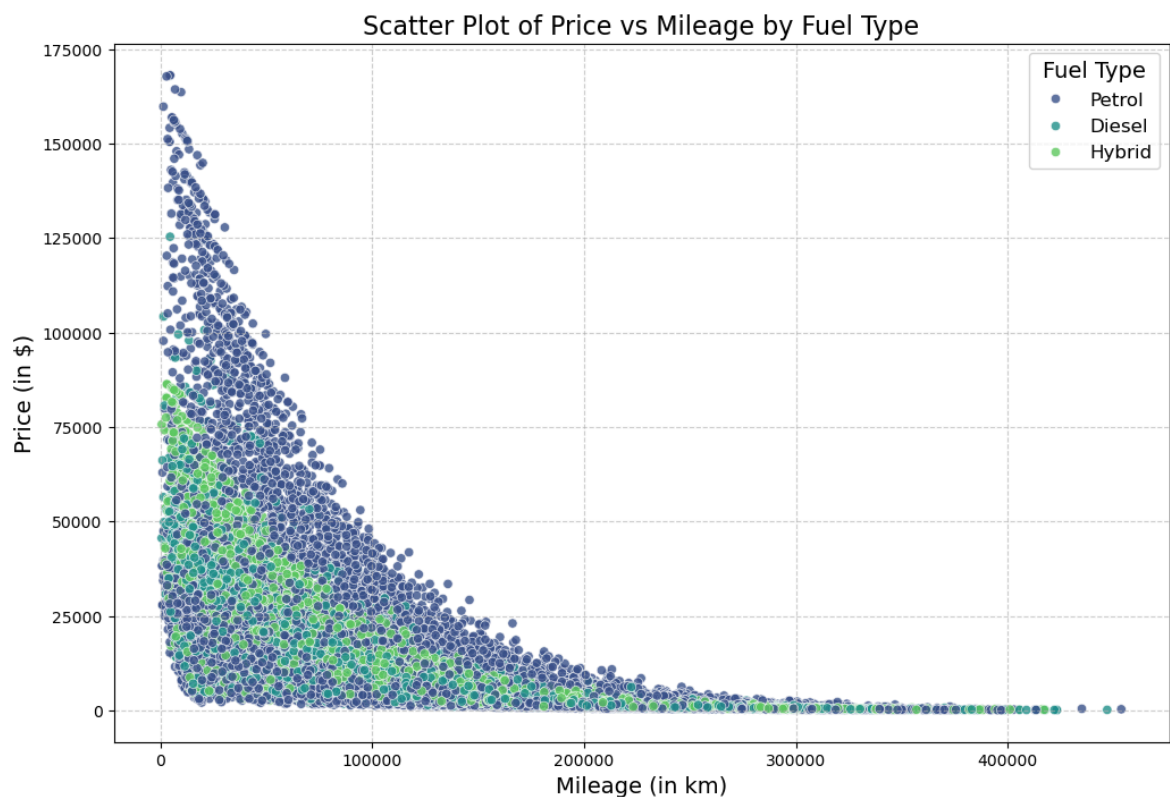
Produce a scatter plot of price vs mileage, coloured by fuel type.

```
In [7]: # Scatter plot of price vs mileage, colored by fuel type
import seaborn as sns
import matplotlib.pyplot as plt

# Create the scatter plot
plt.figure(figsize=(12, 8)) # Increased figure size for better clarity
sns.scatterplot(
    data=car_sales,
    x='Mileage',
    y='Price',
    hue='Fuel type',
    palette='viridis',
    alpha=0.8, # Slightly reduced transparency for better visibility
    edgecolor='w', # Added white edge color for better distinction of points
    linewidth=0.5 # Thin edge lines for clarity
)

# Add titles and labels
plt.title('Scatter Plot of Price vs Mileage by Fuel Type', fontsize=16)
plt.xlabel('Mileage (in km)', fontsize=14)
plt.ylabel('Price (in $)', fontsize=14)
plt.legend(title='Fuel Type', fontsize=12, title_fontsize=14)

# Display the plot
plt.grid(True, linestyle='--', alpha=0.6) # Added grid for better readability
plt.show()
```



Answer-6:

The scatter plot of price vs mileage, colored by fuel type, provides insights into the relationship between these variables and the influence of fuel type on pricing.

Code Explanation:

- **Seaborn Library:** The `seaborn` library is used for creating the scatter plot. It simplifies the process of adding color and other visual elements to the plot.
- **Hue Parameter:** The `hue` parameter in `sns.scatterplot` is used to color the points based on the `Fuel type` column. This allows for easy differentiation between fuel types.
- **Alpha and Edgecolor:** The `alpha` parameter adjusts the transparency of the points, making overlapping points more visible. The `edgecolor=None` ensures a clean look without borders around the points.
- **Titles and Labels:** Titles and axis labels are added to make the plot informative and self-explanatory.

Observations:

- **Price vs Mileage Relationship:** The scatter plot reveals a general trend where higher mileage is associated with lower prices. This is expected, as cars with higher mileage typically have more wear and tear, reducing their value.
- **Fuel Type Influence:** The color coding by fuel type highlights differences in pricing trends among fuel types. For example, hybrid cars may show a different price distribution compared to petrol or diesel cars.
- **Outliers:** The plot helps identify outliers, such as cars with very high prices despite high mileage. These could be luxury or high-performance vehicles.

Conclusion:

The scatter plot is a valuable tool for visualizing the interplay between price, mileage, and fuel type. It provides a clear and intuitive way to identify trends, patterns, and anomalies in the data.

Question 7

Compute price per kilometre.

```
In [8]: # Price per kilometre
car_sales['Price_per_km'] = car_sales['Price'] / car_sales['Mileage']
car_sales[['Price', 'Mileage', 'Price_per_km']].head()
```

Out [8]:

	Price	Mileage	Price_per_km
0	3074	127300	0.024148
1	49704	57850	0.859188
2	24072	39190	0.614238
3	1705	210814	0.008088
4	4101	127869	0.032072

Answer-7:

Purpose:

The `Price_per_km` metric provides a standardized way to compare the cost of cars relative to their mileage. It helps identify cars that offer better value for money in terms of price per kilometre driven.

Calculation:

- **Formula:** $\text{Price_per_km} = \text{Price} / \text{Mileage}$
- This formula divides the price of each car by its mileage to compute the cost per kilometre.

Insights:

- Cars with **lower `Price_per_km` values** are generally more cost-effective, assuming other factors like condition and age are similar.
- **Outliers** in `Price_per_km` may indicate luxury or high-performance vehicles that command a premium price despite high mileage.

Considerations:

- Cars with **very low mileage** may have disproportionately high `Price_per_km` values, as the denominator (mileage) is small.
 - Cars with **zero mileage** (if any) would result in a division by zero error. These cases should be handled separately.
-

Question 8

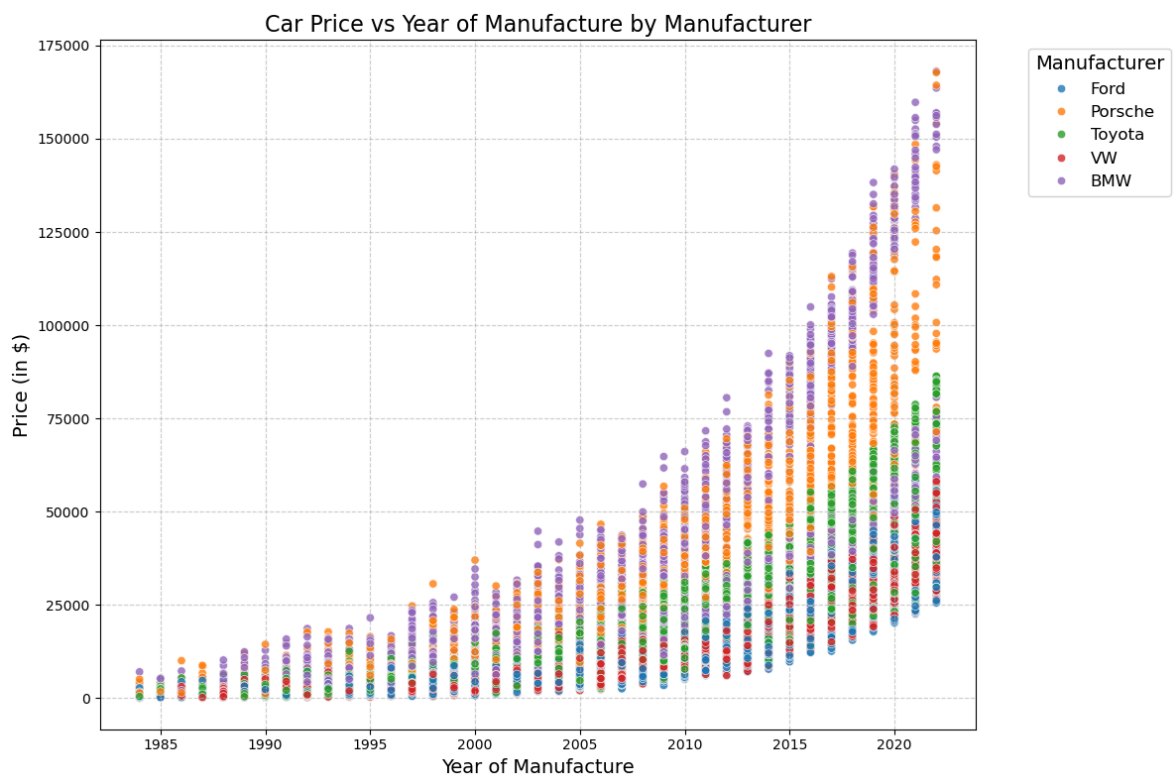
Plot car price vs year of manufacture, colouring points by manufacturer.

```
In [9]: # Price vs year of manufacture
import seaborn as sns
import matplotlib.pyplot as plt

# Create the scatter plot
plt.figure(figsize=(12, 8))
sns.scatterplot(
    data=car_sales,
    x='Year of manufacture',
    y='Price',
    hue='Manufacturer',
    palette='tab10',
    alpha=0.8,
    edgecolor='w',
    linewidth=0.5
)

# Add titles and labels
plt.title('Car Price vs Year of Manufacture by Manufacturer', fontsize=16)
plt.xlabel('Year of Manufacture', fontsize=14)
plt.ylabel('Price (in $)', fontsize=14)
plt.legend(title='Manufacturer', fontsize=12, title_fontsize=14, bbox_to_

# Display the plot
plt.grid(True, linestyle='--', alpha=0.6)
plt.tight_layout()
plt.show()
```



Answer-8:

Code Explanation:

- **Seaborn Library:** The `seaborn` library is used to create the scatter plot, which allows for easy visualization of relationships between variables.
- **Hue Parameter:** The `hue` parameter is used to color the points based on the `Manufacturer` column, enabling differentiation between manufacturers.
- **Alpha and Edgecolor:** Transparency (`alpha=0.8`) and white edge color (`edgecolor='w'`) are applied to improve the clarity of overlapping points.
- **Titles and Labels:** Titles, axis labels, and a legend are added to make the plot informative and visually appealing.

Plot Insights:

- **Price vs Year Relationship:** The plot shows how car prices vary with the year of manufacture. Generally, newer cars tend to have higher prices, reflecting their lower mileage and better condition.
- **Manufacturer Influence:** The color coding highlights differences in pricing trends among manufacturers. For instance, luxury brands like BMW may show consistently higher prices compared to more affordable brands like Toyota.
- **Outliers:** The plot helps identify outliers, such as older cars with unusually high prices. These could be classic or collector cars.

Conclusion:

The scatter plot effectively visualizes the relationship between car price, year of manufacture, and manufacturer. It provides valuable insights into market trends and the influence of brand reputation on pricing.

Question 9

Compare Toyota and BMW.

```
In [10]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Filter data for Toyota and BMW
toyota_bmw = car_sales[car_sales['Manufacturer'].isin(['Toyota', 'BMW'])]

# Summary statistics for Toyota and BMW
summary_stats = toyota_bmw.groupby('Manufacturer')[['Price', 'Mileage', 'Year']]
display(summary_stats) # Use display() to render the table properly

# Average comparison
avg_comparison = toyota_bmw.groupby('Manufacturer')[['Price', 'Mileage', 'Year']]
print("\nAverage Comparison (Mean values):\n")
display(avg_comparison)

# Bar chart for mean values
avg_comparison.plot(kind='bar', figsize=(10, 6), edgecolor='black')
```

```

plt.title("Average Comparison between Toyota and BMW")
plt.ylabel("Mean Value")
plt.xticks(rotation=0)
plt.grid(axis='y', linestyle='--', alpha=0.6)
plt.show()

print(" !! Summary Statistics (All Numeric Columns): !! \n")
display(toyota_bmw.groupby('Manufacturer').describe())

# Fuel Type Distribution

plt.figure(figsize=(7, 5))
sns.countplot(data=toyota_bmw, x='Manufacturer', hue='Fuel type', palette

plt.title("Fuel Type Distribution: Toyota vs BMW")
plt.ylabel("Count")
plt.grid(axis='y', linestyle='--', alpha=0.5)
plt.show()

# Price vs Mileage (Scatter)

plt.figure(figsize=(8, 6))

sns.scatterplot(data=toyota_bmw, x='Mileage', y='Price', hue='Manufacture
plt.title("Price vs Mileage: Toyota vs BMW")
plt.xlabel("Mileage")
plt.ylabel("Price")
plt.grid(True, linestyle='--', alpha=0.5)
plt.show()

# Yearly Average Price Trend

plt.figure(figsize=(8, 6))
sns.lineplot(data=toyota_bmw, x='Year of manufacture', y='Price', hue='Ma

plt.title("Average Price Trend over Years: Toyota vs BMW")
plt.xlabel("Year of Manufacture")
plt.ylabel("Average Price")

plt.grid(True, linestyle='--', alpha=0.5)
plt.show()

# Visual comparison of Price, Mileage, and Engine Size
fig, axes = plt.subplots(1, 3, figsize=(18, 6))

# Price comparison
sns.boxplot(data=toyota_bmw, x='Manufacturer', y='Price', ax=axes[0], hue
axes[0].set_title('Price Comparison')
axes[0].set_ylabel('Price (in $)')

# Mileage comparison
sns.boxplot(data=toyota_bmw, x='Manufacturer', y='Mileage', ax=axes[1], h
axes[1].set_title('Mileage Comparison')
axes[1].set_ylabel('Mileage (in km)')

# Engine size comparison
sns.boxplot(data=toyota_bmw, x='Manufacturer', y='Engine size', ax=axes[2]
axes[2].set_title('Engine Size Comparison')
axes[2].set_ylabel('Engine Size (in L)')

```

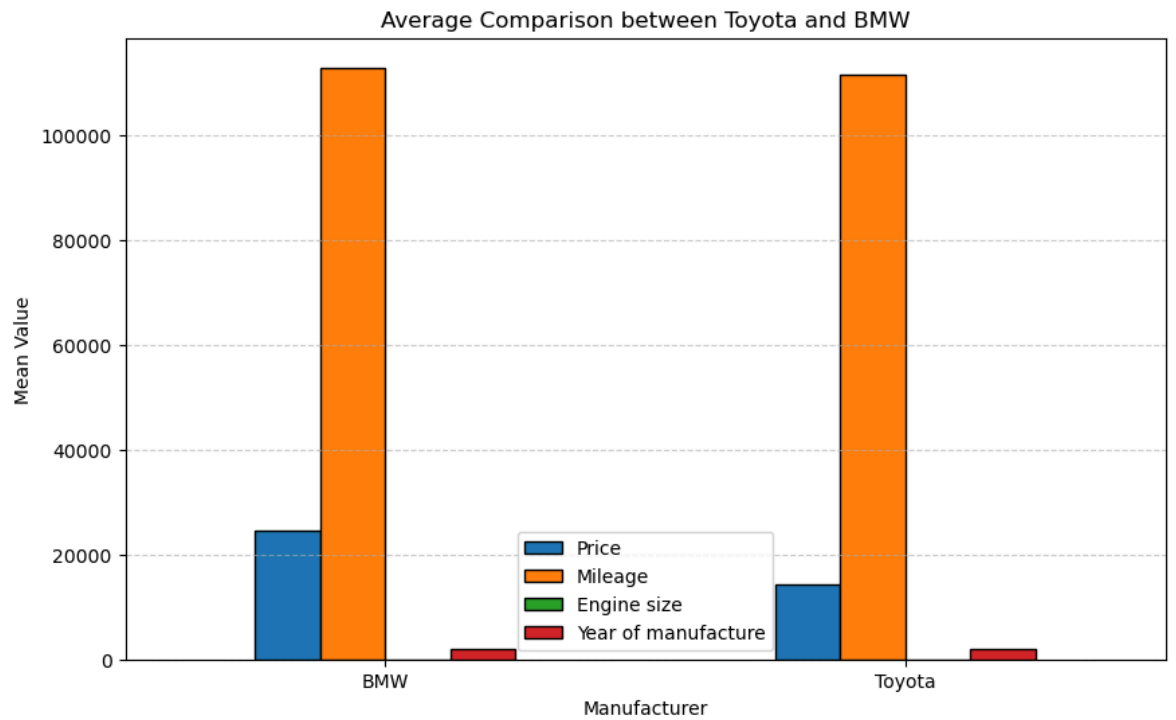
```
plt.tight_layout()
plt.show()
```

	count	mean	std	min	25%	50%	75%
Manufacturer							
BMW	4965.0	24429.459215	27901.217685	167.0	5318.00	14384.0	33592.0
Toyota	12554.0	14340.362275	14706.876784	176.0	3373.25	8802.0	20918.0

2 rows × 24 columns

Average Comparison (Mean values):

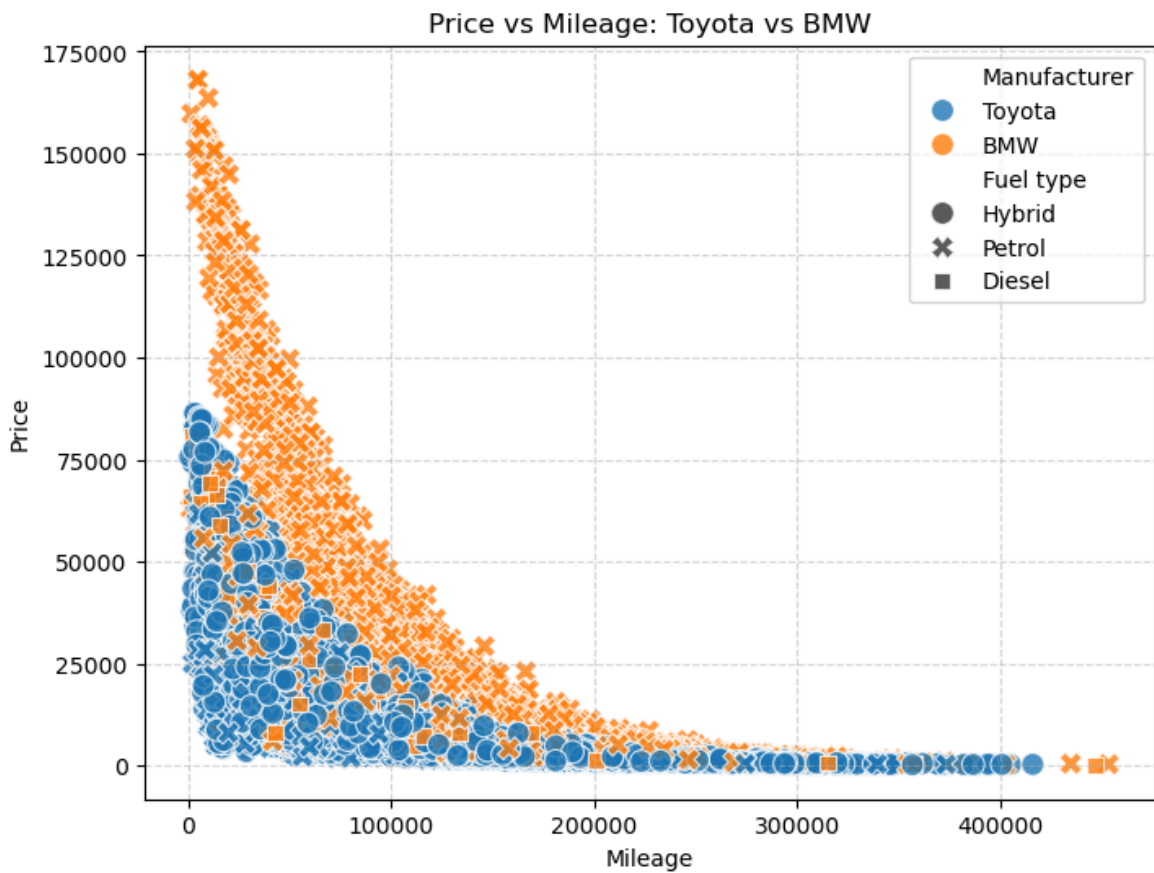
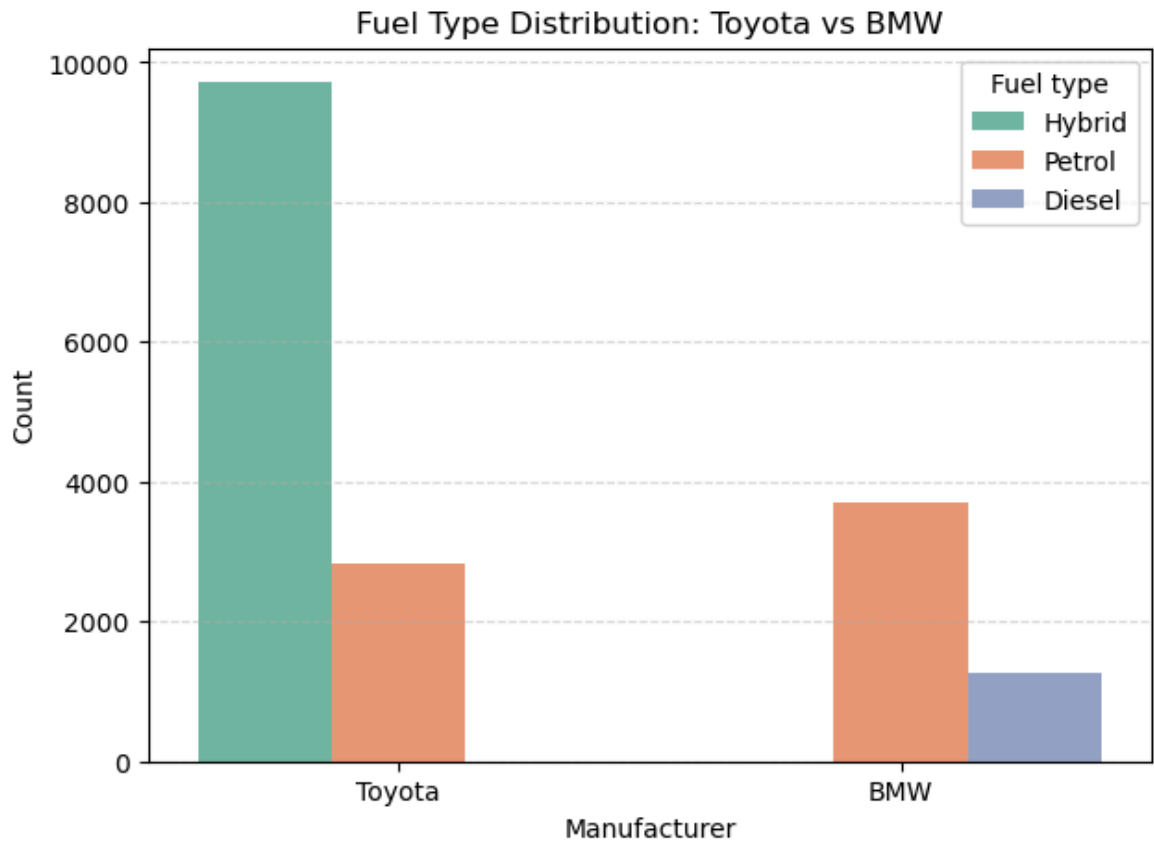
	Price	Mileage	Engine size	Year of manufacture
Manufacturer				
BMW	24429.459215	112837.733333	3.102437	2004.048137
Toyota	14340.362275	111361.126494	1.574367	2004.316951

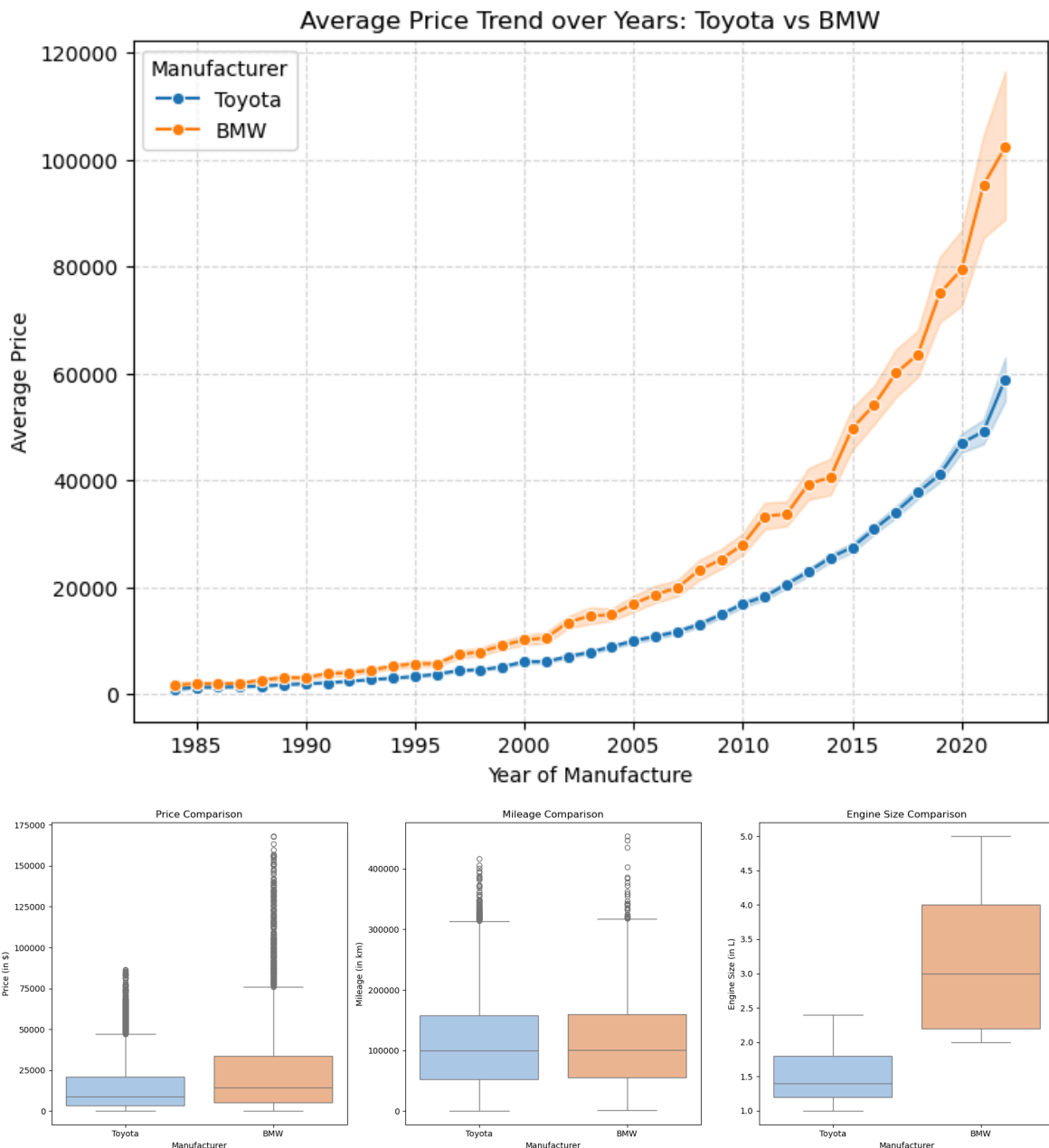


!! Summary Statistics (All Numeric Columns): !!

	count	mean	std	min	25%	50%	75%	max	count	Engine size	Year of m
Manufacturer											
BMW	4965.0	3.102437	1.030934	2.0	2.2	3.0	4.0	5.0	4965.0	2004.048137	
Toyota	12554.0	1.574367	0.456752	1.0	1.2	1.4	1.8	2.4	12554.0	2004.316951	

2 rows × 40 columns





Answer-9:

Code Explanation:

Filtering Data:

The analysis begins by filtering the dataset to include only Toyota and BMW vehicles using the `isin()` method on the `Manufacturer` column. This creates the `toyota_bmw` DataFrame, ensuring that all subsequent analyses focus exclusively on these two manufacturers.

Summary Statistics:

The first `groupby('Manufacturer').describe()` call provides detailed summary statistics for key numerical features (`Price` , `Mileage` , and `Engine size`). This includes measures of central tendency (mean, median) and dispersion (standard deviation, interquartile range), offering a quantitative overview of the differences between Toyota and BMW.

Average Comparison:

The mean values for `Price` , `Mileage` , `Engine size` , and `Year of manufacture` are calculated and visualized using a bar chart. This provides a clear comparison of the average characteristics of Toyota and BMW vehicles.

Overall Summary Statistics:

A broader `groupby('Manufacturer').describe()` call is used to display descriptive statistics for all numeric columns in the filtered dataset. This ensures that no quantifiable attribute is overlooked in the analysis.

Fuel Type Distribution:

The `sns.countplot` function is used to visualize the distribution of the categorical variable `Fuel type` for each manufacturer. This highlights the composition of fuel types (e.g., hybrid, petrol, diesel) within Toyota and BMW vehicles.

Price vs. Mileage:

A scatter plot is generated using `sns.scatterplot` to visualize the relationship between `Mileage` and `Price` . The `hue='Manufacturer'` and `style='Fuel type'` parameters allow for a detailed comparison of how the depreciation curve differs across brands and fuel types.

Yearly Average Price Trend:

The `sns.lineplot` function is used to plot the average `Price` against `Year of manufacture` . This visualizes the depreciation trend over time for each brand, highlighting the dynamic price gap between newer and older models.

Visual Comparison:

Three box plots are created using `sns.boxplot` to compare the distributions of `Price`, `Mileage`, and `Engine size` for Toyota and BMW. These plots highlight the median, interquartile range (IQR), and potential outliers for each feature.

Insights:

Price and Variability:

- BMW cars consistently command a significantly higher average price and exhibit greater variability in price compared to Toyota.
 - This reflects BMW's premium positioning and diverse range of high-spec models.
-

Mileage and Age:

- Toyota cars have higher average mileage and are generally older, aligning with their reputation for long-term reliability and lower cost of ownership.
 - BMW models are typically newer with lower mileage, reflecting their luxury market focus.
-

Engine Size and Performance:

- BMW vehicles have larger average engine sizes, consistent with their emphasis on high-performance driving experiences.
 - Toyota models lean toward smaller, more fuel-efficient engines.
-

Fuel Type Strategy:

- Toyota's strength in the hybrid market is evident, while BMW primarily focuses on petrol and diesel vehicles, showcasing different technological priorities.
-

Depreciation Curve:

- The scatter plot and line plot reveal that BMW maintains its price premium across age and mileage, though both brands show the expected negative correlation (price decreases as mileage/age increases).
-

Conclusion:

This comprehensive analysis highlights the distinct market positions of Toyota and BMW:

- **Toyota:** Focuses on affordability, reliability, and fuel efficiency, making it a popular choice for cost-conscious buyers.
- **BMW:** Targets the luxury and performance segment, characterized by higher prices, newer models, and larger engines.

The multi-faceted comparison provides a clear, data-driven profile of the two brands' relative strengths in the used car market.
