

```
In [11]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings

In [17]: plt.rcParams['figure.figsize'] = [19, 8]
warnings.filterwarnings('ignore')

In [8]: try:
cars_df = pd.read_csv("C:\\Users\\megha\\Downloads\\cars_df.csv", encoding="utf-8", nrows=20)
print("Dataset loaded successfully.")
except FileNotFoundError:
print("Error: 'cars_df.csv' not found. Please ensure the file is in the correct directory and update the path.")
if not cars_df.empty:
print("--- Initial Data Details ---")
print("Dataset shape (rows, columns):", (cars_df.shape))
print("Column names:", cars_df.columns)
cars_df.info()

Dataset loaded successfully.

--- Initial Data Details ---
Dataset shape (rows, columns): (20, 15)

Column Info:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 20 entries, 0 to 19
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   Make        20 non-null      object
 1   Model       20 non-null      object
 2   Type        20 non-null      object
 3   Origin      20 non-null      object
 4   DriveTrain  20 non-null      object
 5   MSRP        20 non-null      float64
 6   Invoice     20 non-null      float64
 7   EngineSize  20 non-null      float64
 8   Cylinders   20 non-null      float64
 9   Horsepower  20 non-null      float64
10  MPG_City    20 non-null      float64
11  MPG_Highway 20 non-null      float64
12  Weight      20 non-null      float64
13  Wheelbase   20 non-null      float64
14  Length      20 non-null      float64
dtypes: float64(10), object(5)
memory usage: 2.5+ MB

In [6]: print("\nFirst 5 rows of the dataset:")
print(cars_df.head())

First 5 rows of the dataset:
   Make      Model      Type      Origin  DriveTrain  MSRP  Invoice  \
0  Acura      MDX        SUV        Asia       All      36945.0  33337.0
1  Acura  RSX Type S 2dr  Sedan      Asia       Front    23820.0  21761.0
2  Acura      TSX 4dr  Sedan      Asia       Front    26990.0  24647.0
3  Acura      TL 4dr  Sedan      Asia       Front    33195.0  30299.0
4  Acura      3.5 RL 4dr  Sedan      Asia       Front    43755.0  39014.0

   EngineSize  Cylinders  Horsepower  MPG_City  MPG_Highway  Weight  \
0         3.5         6.0         265.0        17.0         23.0  4451.0
1         2.0         4.0         200.0        24.0         31.0  2778.0
2         2.4         4.0         200.0        22.0         29.0  3230.0
3         3.2         6.0         270.0        20.0         28.0  3575.0
4         3.5         6.0         225.0        18.0         24.0  3880.0

   Wheelbase  Length
0      106.0   189.0
1      101.0   172.0
2      105.0   183.0
3      108.0   186.0
4      115.0   197.0

In [9]: print("\nLast 5 rows of the dataset:")
print(cars_df.tail())

Last 5 rows of the dataset:
   Make      Model      Type      Origin  DriveTrain  MSRP  Invoice  \
15 Audi  A4 3.0 Quattro convertible 2dr  Sedan      Europe  All    44240.0
16 Audi  A6 2.7 Turbo Quattro 4dr  Sedan      Europe  All    42840.0
17 Audi  A6 4.2 Quattro 4dr  Sedan      Europe  All    49690.0
18 Audi  A8 1 Quattro 4dr  Sedan      Europe  All    69190.0
19 Audi  S4 Quattro 4dr  Sedan      Europe  All    48040.0

   Invoice  EngineSize  Cylinders  Horsepower  MPG_City  MPG_Highway  Weight  \
15  40075.0         3.0         6.0         220.0        18.0         25.0  4013.0
16  38840.0         2.7         6.0         250.0        18.0         25.0  3836.0
17  44936.0         4.2         8.0         300.0        17.0         24.0  4024.0
18  64740.0         4.2         8.0         330.0        17.0         24.0  4399.0
19  43556.0         4.2         8.0         340.0        14.0         20.0  3825.0

   Wheelbase  Length
15      105.0   180.0
16      109.0   192.0
17      109.0   193.0
18      121.0   204.0
19      104.0   179.0

In [10]: print("\n--- Descriptive Statistics for Numerical Columns ---")
print(cars_df.describe())

--- Descriptive Statistics for Numerical Columns ---
count      MSRP      Invoice  EngineSize  Cylinders  Horsepower  \
0      20.000000      20.000000      20.000000      20.000000      20.000000
mean    41748.500000    37817.150000      3.060000      5.900000    218.750000
std     15151.963838    13665.271774      0.703675      1.209611    47.040716
min     23820.000000    21761.000000      1.800000      4.000000    170.000000
25%     33771.250000    30349.250000      2.925000      6.000000    220.000000
50%     38292.500000    34664.500000      3.000000      6.000000    220.000000
75%     44705.000000    40331.250000      3.500000      6.000000    266.250000
max     89765.000000    79976.000000      4.200000      8.000000    340.000000

count      MPG_City  MPG_Highway  Weight  Wheelbase  Length
0      20.000000      20.000000      20.000000      20.000000      20.000000
mean     18.900000     26.000000    3693.700000    107.100000    185.300000
std       2.44734      2.846974     403.228302      5.066921     8.578233
min      14.000000     20.000000    2778.000000    100.000000    172.000000
25%      17.000000     24.000000    3536.250000    104.000000    179.000000
50%      18.000000     25.000000    3728.000000    105.000000    181.500000
75%      20.000000     28.000000    3883.250000    109.000000    192.000000
max      24.000000    31.000000    4451.000000    121.000000    204.000000

In [13]: print("\nNull values before handling:")
cars_df = pd.read_csv("C:\\Users\\megha\\Downloads\\cars_df.csv", na_values=["NA", "Na", "?", "null", " "])
print(cars_df.isnull().sum())

print(cars_df.rename(columns={'MSRP': 'MRP', 'MPG_City': 'Mileage_City',
'MPG_Highway': 'Mileage_Highway'}, inplace=True))
print("Columns renamed. New column names:")
print(cars_df.columns)

print("\n--- Data Cleaning ---")
initial_row = len(cars_df)
print(f"Number of duplicate rows found: {cars_df.duplicated().sum()}")
cars_df.drop_duplicates(inplace=True)
print(f"Duplicate rows dropped. Dataset shape is now: {cars_df.shape}")

Null values before handling:
Make      0
Model     0
Type      0
Origin    0
DriveTrain 0
MSRP      0
Invoice   0
EngineSize 0
Cylinders  2
Horsepower 0
MPG_City   0
MPG_Highway 0
Weight     0
Wheelbase  0
Length     0
dtype: int64

Column renamed. New column names:
Index(['Make', 'Model', 'Type', 'Origin', 'DriveTrain', 'MRP', 'Invoice',
      'EngineSize', 'Cylinders', 'Horsepower', 'Mileage_City',
      'Mileage_Highway', 'Weight', 'Wheelbase', 'Length'],
      dtype='object')

--- Data Cleaning ---
Number of duplicate rows found: 0
Duplicate rows dropped. Dataset shape is now: (428, 15)

In [14]: if 'Cylinders' in cars_df.columns and cars_df['Cylinders'].isnull().any():
median_cylinders = cars_df['Cylinders'].median()
cars_df['Cylinders'].fillna(median_cylinders, inplace=True)
print(f"Missing 'Cylinders' values filled with median value: {median_cylinders}")
print("\nNull values after handling:")
print(cars_df.isnull().sum())

if 'Cylinders' in cars_df.columns:
cars_df['Cylinders'] = cars_df['Cylinders'].astype(np.int64)
print("\n'Cylinders' column data type converted to int64.")
print(cars_df.dtypes)

Missing 'Cylinders' values filled with median value: 6.0

Null values after handling:
Make      0
Model     0
Type      0
Origin    0
DriveTrain 0
MSRP      0
Invoice   0
EngineSize 0
Cylinders  0
Horsepower 0
Mileage_City 0
Mileage_Highway 0
Weight     0
Wheelbase  0
Length     0
dtype: int64

'Cylinders' column data type converted to int64.
Make      object
Model     object
Type      object
Origin    object
DriveTrain object
MRP       float64
Invoice   float64
EngineSize float64
Cylinders int64
Horsepower float64
Mileage_City float64
Mileage_Highway float64
Weight     float64
Wheelbase  float64
Length     float64
dtype: object

In [15]: print("\n--- Generating Visualizations ---")
plt.figure(figsize=(19, 7))
type_counts = cars_df['Type'].value_counts()
plt.title("Number of Cars by Car Type")
plt.bar(type_counts.index, height=type_counts.values)
plt.xlabel("Car Type")
plt.ylabel("No. of Cars")
plt.show()

--- Generating Visualizations ---

Number of Cars by Car Type

250
200
150
100
50
0
Sedan Sports Wagon Truck Hybrid
No. of Cars

In [16]: plt.figure(figsize=(19, 7))
plt.title("Number of Cars by Origin")
sns.countplot(data=cars_df, x="Origin")
plt.show()

Number of Cars by Origin

160
140
120
100
80
60
40
20
0
Asia Europe USA
count

In [17]: plt.figure(figsize=(19, 7))
plt.title("Number of Cars by Drive Train")
sns.countplot(data=cars_df, x="DriveTrain")
plt.show()

Number of Cars by Drive Train

200
150
100
50
0
All Front DriveTrain Rear
count

In [18]: plt.figure(figsize=(19, 7))
plt.title("Frequency Distribution of Cars by Mileage in City")
plt.hist((cars_df['Mileage_City']))
plt.xlabel("Mileage in City")
plt.ylabel("Count of Cars")
plt.show()

Frequency Distribution of Cars by Mileage in City

200
175
150
125
100
75
50
25
0
10 20 30 40 50 60
Count of Cars
Mileage in City

In [19]: plt.figure(figsize=(19, 7))
plt.title("Correlation Matrix of Numerical Features")
numerical_df = cars_df.select_dtypes(include=np.number)
sns.heatmap(numerical_df.corr(), annot=True)
plt.show()

Correlation Matrix of Numerical Features

MRP 1 0.57 0.65 0.83 -0.48 -0.44 0.45 0.15 0.17
Invoice 0.57 1 0.62 0.78 -0.48 -0.44 0.45 0.15 0.17
EngineSize 0.65 0.62 1 0.78 -0.48 -0.44 0.45 0.15 0.17
Cylinders 0.83 0.78 0.78 1 -0.48 -0.44 0.45 0.15 0.17
Horsepower -0.48 -0.48 -0.48 -0.48 1 -0.39 0.45 0.15 0.17
Mileage_City -0.44 -0.44 -0.44 -0.44 -0.39 1 -0.39 0.45 0.17
Mileage_Highway 0.45 0.45 0.45 0.45 -0.39 -0.39 1 -0.39 0.45
Weight 0.15 0.15 0.15 0.15 0.45 0.45 -0.39 1 0.17
Wheelbase 0.17 0.17 0.17 0.17 0.17 0.17 0.45 0.17 1
Length 0.17 0.17 0.17 0.17 0.17 0.17 0.45 0.17 0.17

In [20]: plt.figure(figsize=(19, 7))
plt.title("Boxplot of Numerical Features to Identify Outliers")
sns.boxplot(data=numerical_df)
plt.grid()
plt.show()

Boxplot of Numerical Features to Identify Outliers

200000
175000
150000
125000
100000
75000
50000
25000
0
MRP Invoice EngineSize Cylinders Horsepower Mileage_City Mileage_Highway Weight Wheelbase Length
Outliers

In [21]: print("\n--- Feature Engineering ---")
cars_cat_df = cars_df.select_dtypes(include=np.number)
cars_num_df = cars_df.select_dtypes(include=np.number)
print("Performing One-Hot Encoding on 'Type' column...")
from sklearn.preprocessing import OneHotEncoder
encoder = OneHotEncoder(sparse_output=False)
encoded_type = encoder.fit_transform(cars_cat_df[['Type']])
cars_encoded_df = pd.DataFrame(encoded_type, columns=encoder.get_feature_names_out(['Type']))
print("One-Hot Encoded DataFrame head:")
print(cars_encoded_df.head())

--- Feature Engineering ---
Performing One-Hot Encoding on 'Type' column...
One-Hot Encoded DataFrame head:
Type_Hybrid Type_SUV Type_Sedan Type_Sports Type_Truck Type_Wagon
0 0.0 1.0 0.0 0.0 0.0 0.0
1 0.0 0.0 1.0 0.0 0.0 0.0
2 0.0 0.0 0.0 1.0 0.0 0.0
3 0.0 0.0 0.0 1.0 0.0 0.0
4 0.0 0.0 1.0 0.0 0.0 0.0

In [22]: print("\n--- Data Scaling ---")
print("Applying StandardScaler to numerical data...")
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
cars_num_scaled_df = pd.DataFrame(scaler.fit_transform(cars_num_df), columns=cars_num_df.columns)
print("Scaled numerical data description:")
print(cars_num_scaled_df.describe())

--- Data Scaling ---
Applying StandardScaler to numerical data...
Scaled numerical data description:
count      MRP      Invoice  EngineSize  Cylinders  Horsepower  \
0      4.280000e+02      4.280000e+02      4.280000e+02      4.280000e+02      4.280000e+02
mean    -1.577139e-16      0.000000      3.652322e-16      2.158191e-16      -7.470606e-17
std     1.001170e+00      1.001170e+00      1.001170e+00      1.001170e+00      1.001170e+00
min    -1.128919e+00      -1.142905      -1.712933e+00      -1.408347e+00      -1.991379e+00
25%    -6.409709e-01      -0.632676      -7.421020e-01      -1.164434e+00      -7.091854e-01
50%    -2.027386e-01      -0.267866      -1.776656e-01      1.233648e-01      8.202571e-02
75%      3.312970e-01      0.323216      6.351230e-01      1.233648e-01      5.451340e-01
max      8.227635e+00      8.146034      4.608757e+00      3.986788e+00      3.959670e+00

count      Mileage_City  Mileage_Highway  Weight  Wheelbase  Length
0      4.280000e+02      4.280000e+02      4.280000e+02      4.280000e+02      4.280000e+02
mean    -2.448716e-16      -3.527811e-17      9.130806e-17      6.474572e-16      2.573227e-16
std     1.001170e+00      1.001170e+00      1.001170e+00      1.001170e+00      1.001170e+00
min    -1.128919e+00      -1.142905      -1.712933e+00      -1.408347e+00      -1.991379e+00
25%    -6.409709e-01      -0.632676      -7.421020e-01      -1.164434e+00      -7.091854e-01
50%    -2.027386e-01      -0.267866      -1.776656e-01      1.233648e-01      8.202571e-02
75%      3.312970e-01      0.323216      6.351230e-01      1.233648e-01      5.451340e-01
max      8.227635e+00      8.146034      4.608757e+00      3.986788e+00      3.959670e+00

In [23]: cars_encoded_df.reset_index(drop=True, inplace=True)
cars_num_scaled_df.reset_index(drop=True, inplace=True)
final_df = pd.concat([cars_encoded_df, cars_num_scaled_df], axis=1)
print("\n--- Final Processed DataFrame ---")
print(final_df.shape)
print(final_df.head())

--- Final Processed DataFrame ---
Final shape: (428, 16)
Type_Hybrid Type_SUV Type_Sedan Type_Sports Type_Truck Type_Wagon \
0 0.0 1.0 0.0 0.0 0.0 0.0
1 0.0 0.0 1.0 0.0 0.0 0.0
2 0.0 0.0 0.0 1.0 0.0 0.0
3 0.0 0.0 0.0 1.0 0.0 0.0
4 0.0 0.0 1.0 0.0 0.0 0.0

MRP Invoice EngineSize Cylinders Horsepower Mileage_City \
0 0.214856 0.188537 0.273884 0.123365 0.684503 -0.584995
1 -0.461376 -0.468388 -1.080764 -1.164443 -0.223395 0.752902
2 -0.298050 -0.304611 -0.719525 -1.164443 -0.223395 0.370645
3 0.216147 0.016134 0.002994 0.123365 0.754187 -0.016111
4 0.365724 0.110700 0.287884 0.123365 0.125928 -0.393967

Mileage_Highway Weight Wheelbase Length
0 -0.702319 1.151621 -0.229477 0.183935
1 0.724832 -1.055214 -0.861735 -1.001460
2 0.376065 -0.458983 -0.379929 -0.234440
3 0.201482 -0.031896 -0.018574 -0.025232
4 -0.495852 0.398428 -0.824586 0.743768

In [24]: print("\n--- Outlier Treatment ---")
Q1 = final_df.quantile(0.25)
Q3 = final_df.quantile(0.75)
IQR = Q3 - Q1
lower_limit = Q1 - 1.5 * IQR
upper_limit = Q3 + 1.5 * IQR
rows_before = final_df.shape[0]
outliers_mask = (final_df < lower_limit) | (final_df > upper_limit)
cars_df_no_outliers = final_df[~outliers_mask]
rows_after = cars_df_no_outliers.shape[0]
print(f"Number of outliers detected and removed: {rows_before - rows_after}")
print(f"Shape of DataFrame after removing outliers: {cars_df_no_outliers.shape}")

--- Outlier Treatment ---
Number of outliers detected and removed: 203
Shape of DataFrame after removing outliers: (225, 16)

In [25]: plt.figure(figsize=(19, 7))
sns.boxplot(data=cars_df_no_outliers)
plt.title("Boxplot of Final Data After Outlier Removal")
plt.grid()
plt.show()

Boxplot of Final Data After Outlier Removal

2
1
0
-1
Type_Hybrid Type_SUV Type_Sedan Type_Sports Type_Truck Type_Wagon MRP Invoice EngineSize Cylinders Horsepower Mileage_City Mileage_Highway Weight Wheelbase Length
Analysis complete.
```