

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.ensemble import RandomForestRegressor # Using Regression Model
8 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
9

```

```

1 file_path = "/content/car_price_dataset.csv" # File Path
2 df = pd.read_csv(file_path)
3
4 # Display first 5 rows
5 print("First 5 rows of the dataset:")
6 print(df.head())
7

```

First 5 rows of the dataset:

	Brand	Model	Year	Engine_Size	Fuel_Type	Transmission	Mileage
0	Kia	Rio	2020	4.2	Diesel	Manual	289944
1	Chevrolet	Malibu	2012	2.0	Hybrid	Automatic	5356
2	Mercedes	GLA	2020	4.2	Diesel	Automatic	231440
3	Audi	Q5	2023	2.0	Electric	Manual	160971
4	Volkswagen	Golf	2003	2.6	Hybrid	Semi-Automatic	286618

	Doors	Owner_Count	Price
0	3	5	8501
1	2	3	12092
2	4	2	11171
3	2	1	11780
4	3	3	2867

```

1 # Check Dataset Information
2 print("\nDataset Info:")
3 print(df.info())
4
5 print("\nChecking for Missing Values:")
6 print(df.isnull().sum())
7
8
9 print("\nDataset Statistics:")
10 print(df.describe())
11

```

Dataset Info:

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Brand                 10000 non-null object
1   Model                 10000 non-null object
2   Year                  10000 non-null int64
3   Engine_Size           10000 non-null float64
4   Fuel_Type             10000 non-null object
5   Transmission          10000 non-null object
6   Mileage               10000 non-null int64
7   Doors                 10000 non-null int64
8   Owner_Count           10000 non-null int64
9   Price                 10000 non-null int64
dtypes: float64(1), int64(5), object(4)
memory usage: 781.4+ KB
None

```

Checking for Missing Values:

```

Brand      0
Model      0
Year       0
Engine_Size 0
Fuel_Type  0
Transmission 0
Mileage    0
Doors      0
Owner_Count 0
Price      0
dtype: int64

```

## Dataset Statistics:

	Year	Engine_Size	Mileage	Doors	Owner_Count	\
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	
mean	2011.543700	3.000560	149239.111800	3.497100	2.991100	
std	6.897699	1.149324	86322.348957	1.110097	1.422682	
min	2000.000000	1.000000	25.000000	2.000000	1.000000	
25%	2006.000000	2.000000	74649.250000	3.000000	2.000000	
50%	2012.000000	3.000000	149587.000000	3.000000	3.000000	
75%	2017.000000	4.000000	223577.500000	4.000000	4.000000	
max	2023.000000	5.000000	299947.000000	5.000000	5.000000	

## Price

count	10000.000000
mean	8852.96440
std	3112.59681
min	2000.000000
25%	6646.000000
50%	8858.500000
75%	11086.500000
max	18301.000000

```

1 # Separate numeric and categorical columns
2 numeric_cols = df.select_dtypes(include=['number']).columns
3 categorical_cols = df.select_dtypes(include=['object']).columns
4
5 # Fill missing values for numeric columns with mean
6 df[numeric_cols] = df[numeric_cols].fillna(df[numeric_cols].mean())
7
8 # Fill missing values for categorical columns with the most frequent value (mode)
9 df[categorical_cols] = df[categorical_cols].fillna(df[categorical_cols].mode().iloc[0])
10
11 # Check if missing values are handled
12 print("\nMissing values after filling:")
13 print(df.isnull().sum())
14

```



## Missing values after filling:

Brand	0
Model	0
Year	0
Engine_Size	0
Fuel_Type	0
Transmission	0
Mileage	0
Doors	0
Owner_Count	0
Price	0
dtype:	int64

```

1 df = pd.get_dummies(df, drop_first=True) # Convert categorical variables into numeric
2 X = df.drop(columns=['Price']) # Assuming 'price' is the target variable
3 y = df['Price']
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
6 scaler = StandardScaler()
7 X_train = scaler.fit_transform(X_train)
8 X_test = scaler.transform(X_test)
9
10 model = RandomForestRegressor(n_estimators=100, random_state=42)
11 model.fit(X_train, y_train)
12
13 y_pred = model.predict(X_test)
14

```

```

1 mae = mean_absolute_error(y_test, y_pred)
2 mse = mean_squared_error(y_test, y_pred)
3 rmse = np.sqrt(mse)
4 r2 = r2_score(y_test, y_pred)
5
6 print("\nModel Performance Metrics:")
7 print(f"Mean Absolute Error (MAE): {mae}")
8 print(f"Mean Squared Error (MSE): {mse}")
9 print(f"Root Mean Squared Error (RMSE): {rmse}")
10 print(f"R² Score: {r2}")
11

```



## Model Performance Metrics:

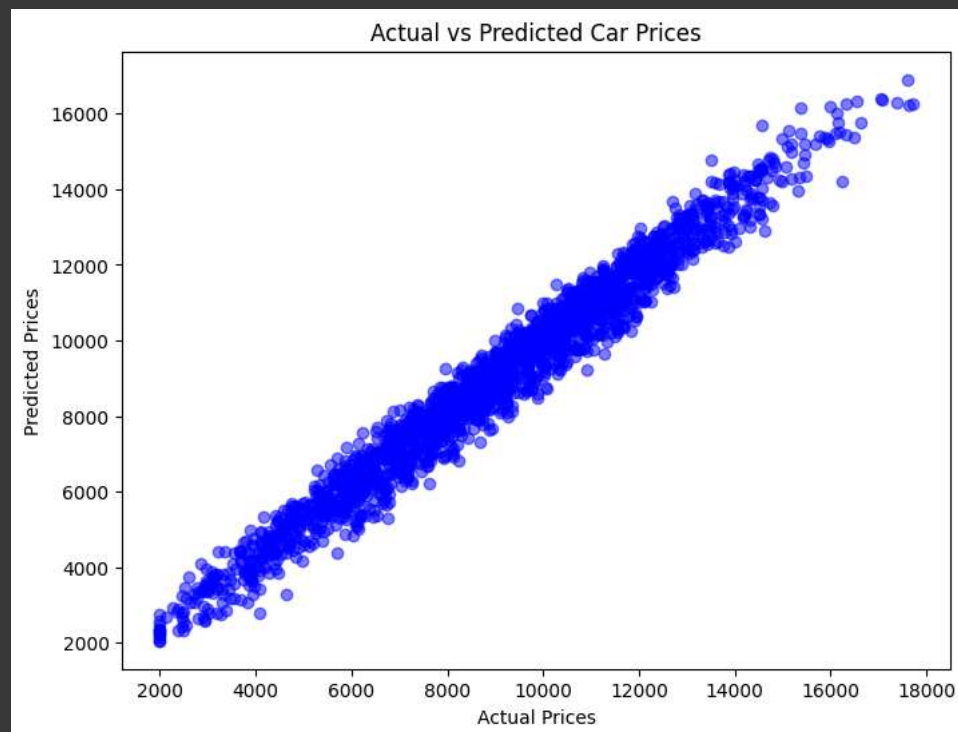
Mean Absolute Error (MAE): 439.79915

Mean Squared Error (MSE): 296942.1043499

Root Mean Squared Error (RMSE): 544.9239436379172

R<sup>2</sup> Score: 0.9676808679542888

```
1 plt.figure(figsize=(8,6))
2 plt.scatter(y_test, y_pred, alpha=0.5, color='blue')
3 plt.xlabel("Actual Prices")
4 plt.ylabel("Predicted Prices")
5 plt.title("Actual vs Predicted Car Prices")
6 plt.show()
7
```



```
1 import joblib
2
3 # Save Model
4 joblib.dump(model, 'car_price_model.pkl')
5
6 # Load Model
7 loaded_model = joblib.load('car_price_model.pkl')
8
```