

# Car Price Prediction and Deployment Using Streamlit Cloud

## Introduction

This project aims to develop a machine learning model to predict car prices based on various features such as brand, model, year, engine size, fuel type, transmission, mileage, doors, and owner count. The model will be deployed using streamlit cloud integrated with Github for interactive usage and demonstration.

## Primary Objective

To build an accurate car price prediction model and deploy it using streamlit, creating an end-to-end data science solution that can be shared and demonstrated to stakeholders.

In [ ]:

## Import Necessary Libraries

In [1]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
import warnings
warnings.filterwarnings('ignore')
```

## Load and Explore the Data

In [2]:

```
df = pd.read_csv(r"C:\Users\DELL\Desktop\Data Analytics Projects\Car Price Prediction\car_price_dataset.csv")
```

In [3]:

```
df.head()
```

Out[3]:

	Brand	Model	Year	Engine_Size	Fuel_Type	Transmission	Mileage	Doors	Owner_Count	Price
0	Kia	Rio	2020	4.2	Diesel	Manual	289944	3	5	8501
1	Chevrolet	Malibu	2012	2.0	Hybrid	Automatic	5356	2	3	12092
2	Mercedes	GLA	2020	4.2	Diesel	Automatic	231440	4	2	11171
3	Audi	Q5	2023	2.0	Electric	Manual	160971	2	1	11780
4	Volkswagen	Golf	2003	2.6	Hybrid	Semi-Automatic	286618	3	3	2867

In [ ]:

In [4]:

```
# Display basic information about the dataset
print("Dataset Shape:", df.shape)
print()
print("\nDataset Info:")
print()
print(df.info())
```

Dataset Shape: (10000, 10)

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
```

In [ ]:

## Data Cleaning and Preprocessing

In [5]: *# Check for missing values*

```
print("Missing values:")
print(df.isnull().sum())
```

```
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
```

In [ ]:

In [6]: *# Check for duplicates*

```
print(f"\nDuplicate rows: {df.duplicated().sum()}")
df = df.drop_duplicates()
```

Duplicate rows: 0

In [ ]:

In [7]: *# Check data types and convert if necessary*

```
print("\nData types:")
print(df.dtypes)
```

```
Data types:
Brand            object
Model            object
Year             int64
Engine_Size      float64
Fuel_Type        object
Transmission     object
Mileage          int64
Doors            int64
Owner_Count      int64
Price            int64
dtype: object
```

In [ ]:

In [8]: *# Check for zero or negative values where inappropriate*

```
print(f"\nCars with zero or negative price: {(df['Price'] <= 0).sum()}")
print(f"Cars with negative mileage: {(df['Mileage'] < 0).sum()}")
print(f"Cars with zero doors: {(df['Doors'] == 0).sum()}")
```

Cars with zero or negative price: 0  
Cars with negative mileage: 0  
Cars with zero doors: 0

In [ ]:

## Exploratory Data Analysis (EDA)

```
In [9]: # Statistical summary
print("Statistical Summary:")
print(df.describe())
```

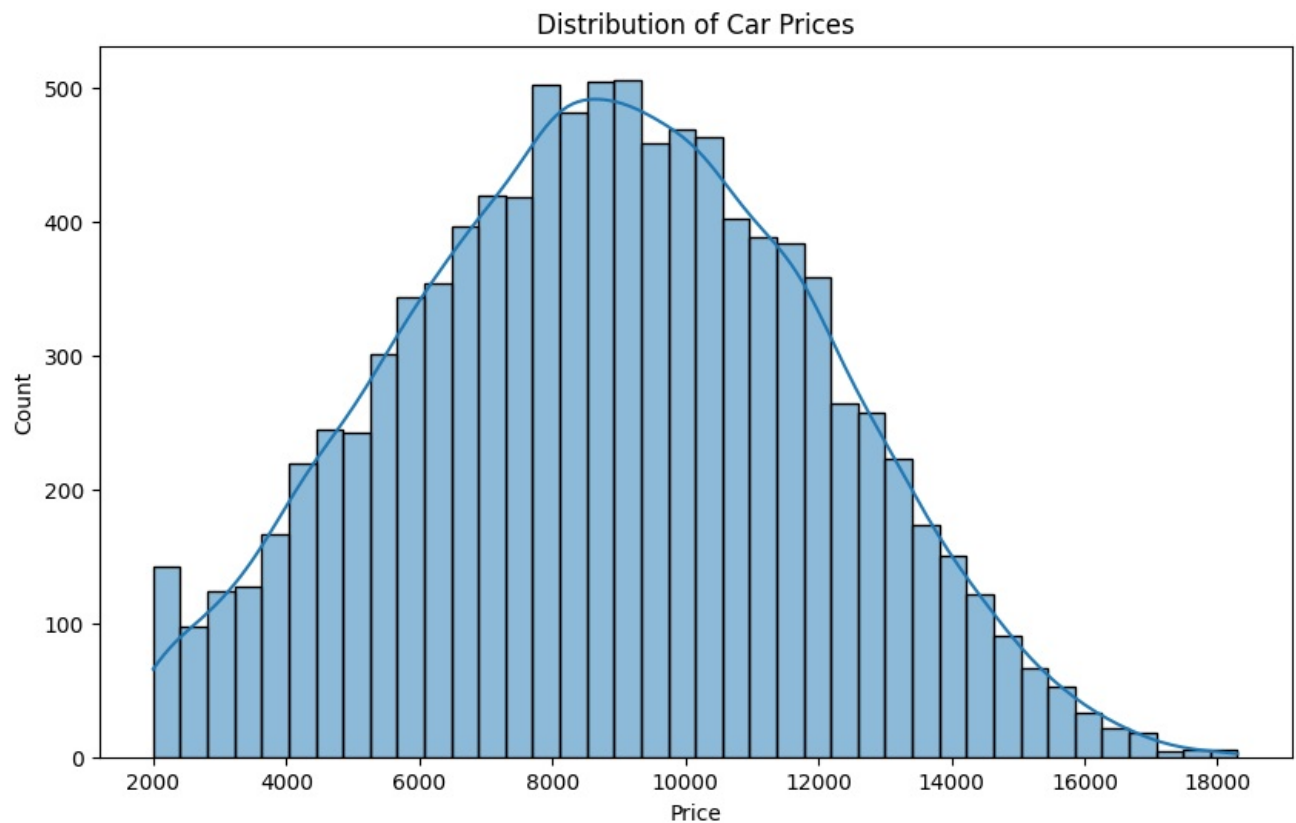
Statistical Summary:

	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

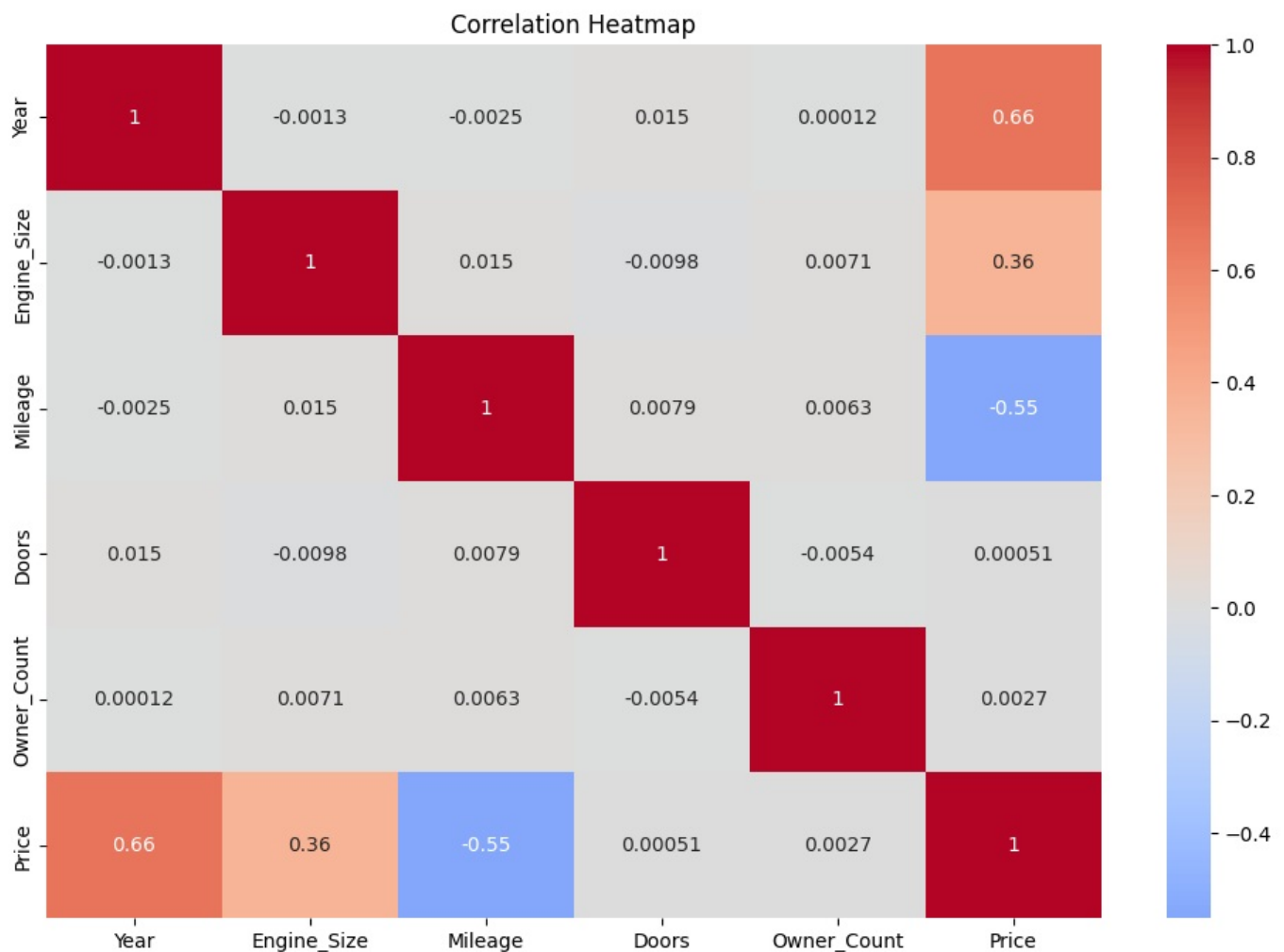
In [ ]:

```
In [10]: # Distribution of target variable
plt.figure(figsize=(10, 6))
sns.histplot(df['Price'], kde=True)
plt.title('Distribution of Car Prices')
plt.show()
```



In [ ]:

```
In [11]: # Correlation heatmap
plt.figure(figsize=(12, 8))
numeric_df = df.select_dtypes(include=[np.number])
sns.heatmap(numeric_df.corr(), annot=True, cmap='coolwarm', center=0)
plt.title('Correlation Heatmap')
plt.show()
```



Based on the correlation heatmap:

- Production Year shows the strongest positive correlation with Price (0.667), indicating that newer cars generally have higher prices.
- Mileage has a moderate negative correlation with Price (-0.55), suggesting that higher mileage typically leads to lower car prices, which aligns with real-world expectations.

In [ ]:

In [ ]:

In [12]: *# Categorical variables analysis*

```

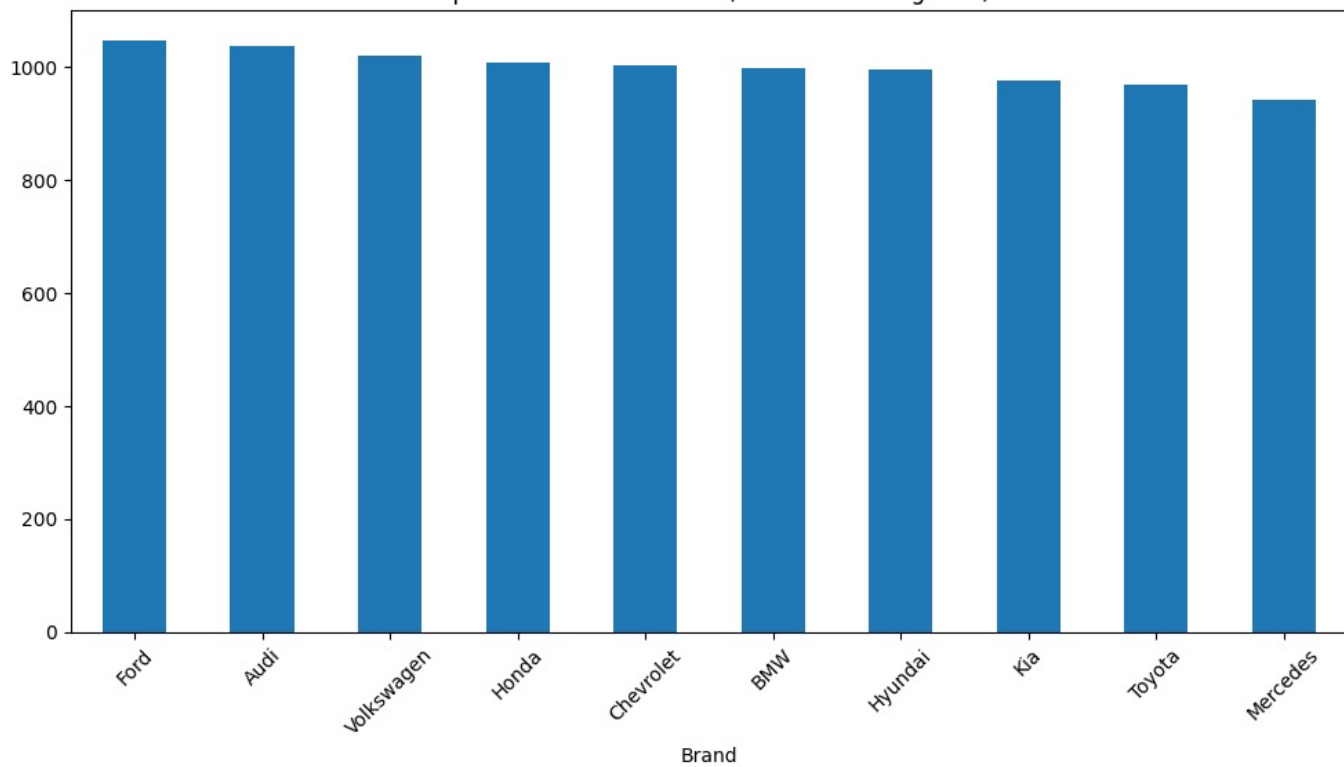
categorical_cols = ['Brand', 'Model', 'Fuel_Type', 'Transmission']
for col in categorical_cols:
    plt.figure(figsize=(10, 6))
    value_counts = df[col].value_counts()
    unique_count = len(value_counts)

    # Dynamic title based on category count
    if unique_count <= 8:
        plot_data = value_counts
        title = f'{col} Distribution ({unique_count} categories)'
    else:
        plot_data = value_counts.head(10)
        title = f'Top 10 {col} Distribution (of {unique_count} total categories)'

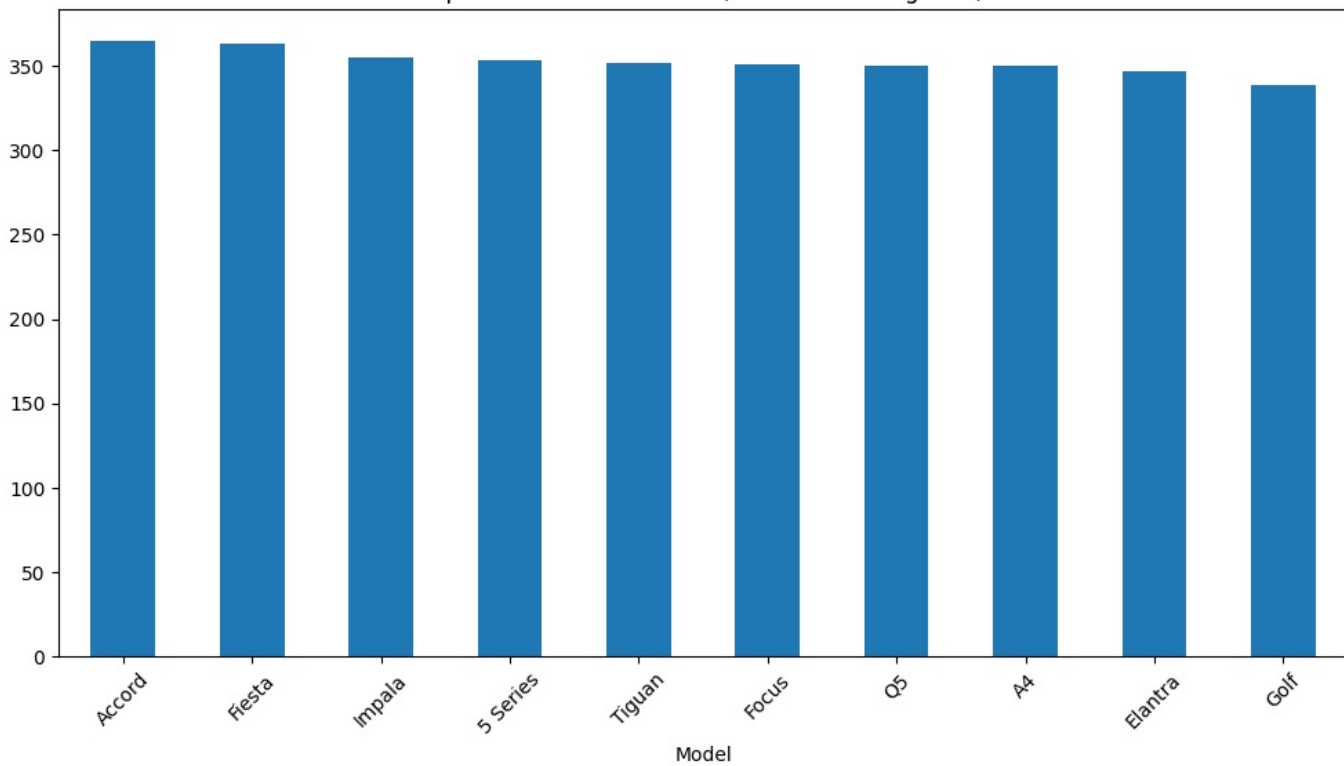
    plot_data.plot(kind='bar')
    plt.title(title)
    plt.xticks(rotation=45)
    plt.tight_layout()
    plt.show()

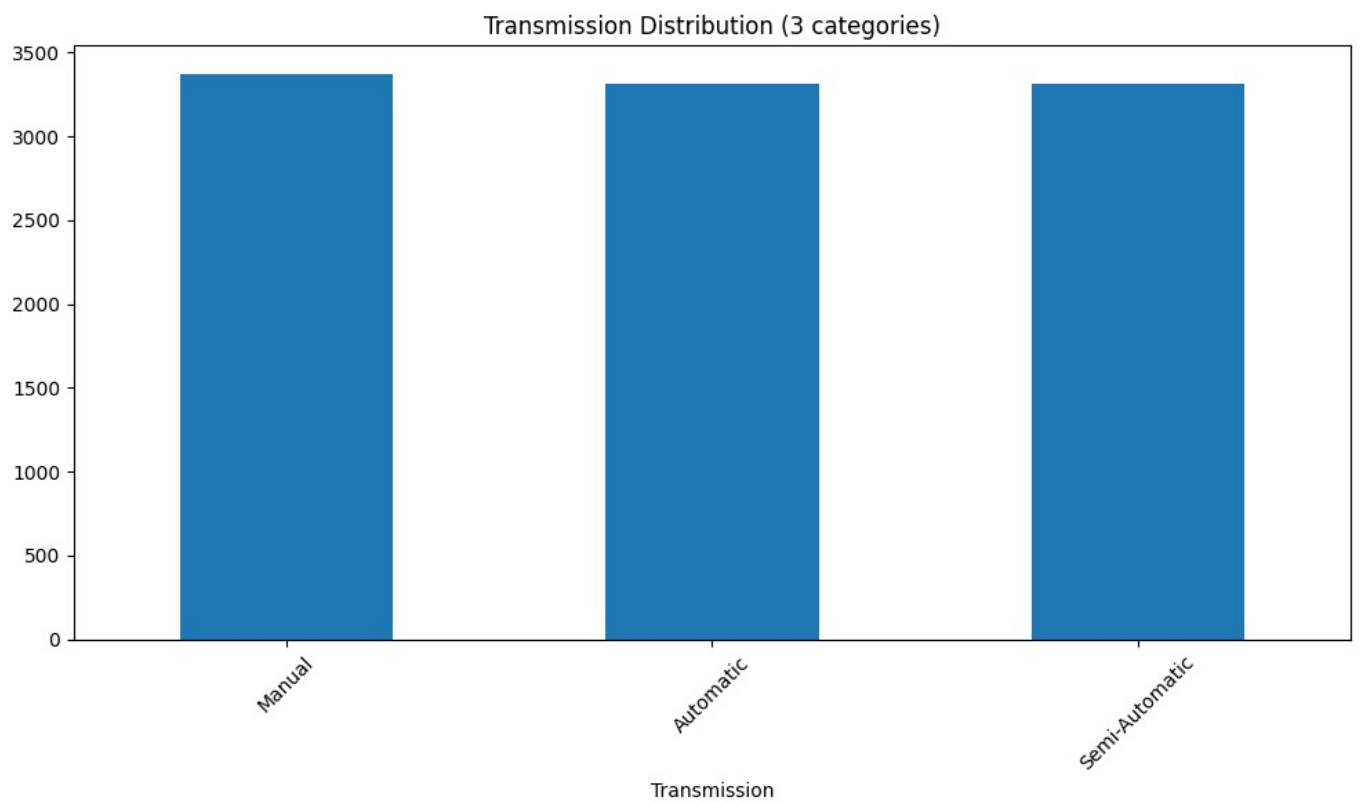
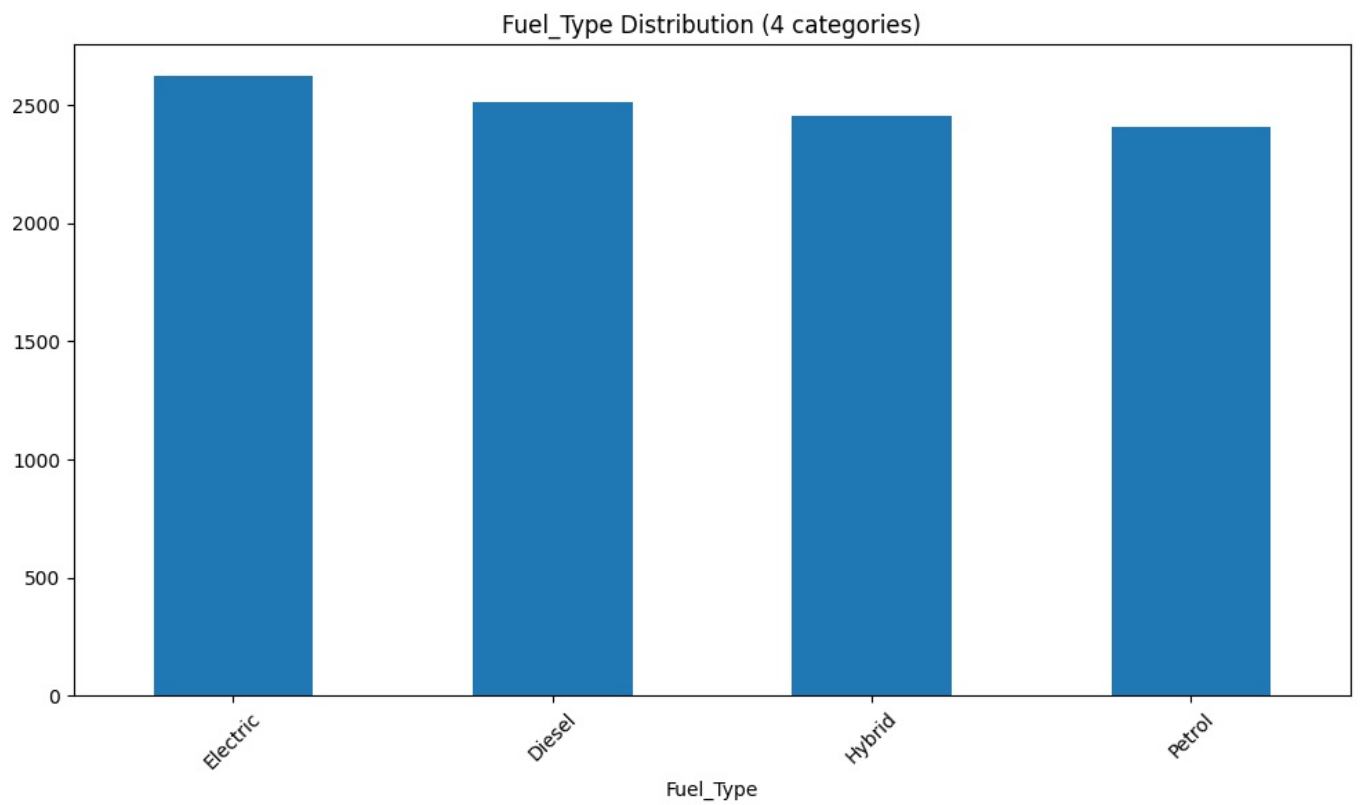
```

Top 10 Brand Distribution (of 10 total categories)



Top 10 Model Distribution (of 30 total categories)





In [ ]:

## Feature Engineering

In [13]: *# Create age feature from year*

```
df['Car_Age'] = 2024 - df['Year'] # Update year as needed
```

In [14]: `df.head()`

Out[14]:

	Brand	Model	Year	Engine_Size	Fuel_Type	Transmission	Mileage	Doors	Owner_Count	Price	Car_Age
0	Kia	Rio	2020	4.2	Diesel	Manual	289944	3	5	8501	4
1	Chevrolet	Malibu	2012	2.0	Hybrid	Automatic	5356	2	3	12092	12
2	Mercedes	GLA	2020	4.2	Diesel	Automatic	231440	4	2	11171	4
3	Audi	Q5	2023	2.0	Electric	Manual	160971	2	1	11780	1
4	Volkswagen	Golf	2003	2.6	Hybrid	Semi-Automatic	286618	3	3	2867	21

In [ ]:

In [15]:

df.head()

Out[15]:

	Brand	Model	Year	Engine_Size	Fuel_Type	Transmission	Mileage	Doors	Owner_Count	Price	Car_Age
0	Kia	Rio	2020	4.2	Diesel	Manual	289944	3	5	8501	4
1	Chevrolet	Malibu	2012	2.0	Hybrid	Automatic	5356	2	3	12092	12
2	Mercedes	GLA	2020	4.2	Diesel	Automatic	231440	4	2	11171	4
3	Audi	Q5	2023	2.0	Electric	Manual	160971	2	1	11780	1
4	Volkswagen	Golf	2003	2.6	Hybrid	Semi-Automatic	286618	3	3	2867	21

In [ ]:

In [16]:

```
# Manual target encoding with comprehensive error handling
categorical_features = ['Brand', 'Model', 'Fuel_Type', 'Transmission']

# Verify all columns exist
available_columns = [col for col in categorical_features if col in df.columns]
print(f"Available categorical columns: {available_columns}")

if len(available_columns) != len(categorical_features):
    missing = set(categorical_features) - set(available_columns)
    print(f"Warning: Missing columns {missing}")

for feature in available_columns:
    try:
        print(f"\nProcessing {feature}...")

        # Calculate mean price for each category
        mean_prices = df.groupby(feature)['Price'].mean().to_dict()
        print(f"Found {len(mean_prices)} unique categories in {feature}")

        # Map the mean prices to create encoded column
        df[f'{feature}_encoded'] = df[feature].map(mean_prices)

        # Handle any NaN values (categories not in training)
        nan_count = df[f'{feature}_encoded'].isna().sum()
        if nan_count > 0:
            overall_mean = df['Price'].mean()
            df[f'{feature}_encoded'].fillna(overall_mean, inplace=True)
            print(f"Filled {nan_count} NaN values with overall mean price: ${overall_mean:.2f}")

        print(f"Successfully encoded {feature}")

    except Exception as e:
        print(f"Error encoding {feature}: {e}")
        continue

# Update features list
features = [f'{col}_encoded' for col in available_columns] + [
    'Year', 'Engine_Size', 'Mileage', 'Doors', 'Owner_Count', 'Car_Age'
]

print(f"\nFinal features to use: {features}")
print(f"Dataset shape: {df[features].shape}")
```

Available categorical columns: ['Brand', 'Model', 'Fuel\_Type', 'Transmission']

Processing Brand...

Found 10 unique categories in Brand

Successfully encoded Brand

Processing Model...

Found 30 unique categories in Model

Successfully encoded Model

Processing Fuel\_Type...

Found 4 unique categories in Fuel\_Type

Successfully encoded Fuel\_Type

Processing Transmission...

Found 3 unique categories in Transmission

Successfully encoded Transmission

Final features to use: ['Brand\_encoded', 'Model\_encoded', 'Fuel\_Type\_encoded', 'Transmission\_encoded', 'Year', 'Engine\_Size', 'Mileage', 'Doors', 'Owner\_Count', 'Car\_Age']

Dataset shape: (10000, 10)

In [17]: df.head()

Out[17]:

	Brand	Model	Year	Engine_Size	Fuel_Type	Transmission	Mileage	Doors	Owner_Count	Price	Car_Age	Brand_encoded
0	Kia	Rio	2020	4.2	Diesel	Manual	289944	3	5	8501	4	8880.0860
1	Chevrolet	Malibu	2012	2.0	Hybrid	Automatic	5356	2	3	12092	12	9015.6839
2	Mercedes	GLA	2020	4.2	Diesel	Automatic	231440	4	2	11171	4	8980.0870
3	Audi	Q5	2023	2.0	Electric	Manual	160971	2	1	11780	1	8929.3737
4	Volkswagen	Golf	2003	2.6	Hybrid	Semi-Automatic	286618	3	3	2867	21	8928.3774

In [ ]:

In [18]:

```
# Let's verify the encoding makes sense
print("=== VERIFICATION ===")

# Check if Brand encoding matches actual average prices
brand_verification = df.groupby('Brand').agg({
    'Price': 'mean',
    'Brand_encoded': 'first'
}).round(2)

print("Brand Encoding Verification:")
print(brand_verification.head())

# Check Fuel Type encoding
fuel_verification = df.groupby('Fuel_Type').agg({
    'Price': 'mean',
    'Fuel_Type_encoded': 'first'
}).round(2)

print("\nFuel Type Encoding Verification:")
print(fuel_verification)
```

=== VERIFICATION ===

Brand Encoding Verification:

	Price	Brand_encoded
Brand		
Audi	8929.37	8929.37
BMW	8704.07	8704.07
Chevrolet	9015.68	9015.68
Ford	8852.57	8852.57
Honda	8665.60	8665.60

Fuel Type Encoding Verification:

	Price	Fuel_Type_encoded
Fuel_Type		
Diesel	8117.34	8117.34
Electric	10032.22	10032.22
Hybrid	9113.03	9113.03
Petrol	8070.56	8070.56

In [ ]:

In [19]:

```
features = [
    'Brand_encoded', 'Model_encoded', 'Fuel_Type_encoded', 'Transmission_encoded', 'Engine_Size', 'Mileage', 'Doors', 'Owner_Count', 'Car_Age'
]
```



```
X = df[features]
y = df['Price']

print(f"Features shape: {X.shape}")
print("\nFeature ranges:")
print(X.describe())
```

Features shape: (10000, 9)

Feature ranges:

	Brand_encoded	Model_encoded	Fuel_Type_encoded	Transmission_encoded	\
count	10000.000000	10000.000000	10000.000000	10000.000000	
mean	8852.964400	8852.964400	8852.964400	8852.964400	
std	109.728109	146.225550	815.577636	765.708083	
min	8665.596630	8517.327381	8070.561826	8264.266385	
25%	8778.279397	8743.761644	8117.336385	8264.266385	
50%	8852.570611	8867.635783	9113.030167	8363.426157	
75%	8929.373796	8967.330218	10032.220190	9938.252939	
max	9015.683948	9156.320635	10032.220190	9938.252939	

	Engine_Size	Mileage	Doors	Owner_Count	Car_Age
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	3.000560	149239.111800	3.497100	2.991100	12.456300
std	1.149324	86322.348957	1.110097	1.422682	6.897699
min	1.000000	25.000000	2.000000	1.000000	1.000000
25%	2.000000	74649.250000	3.000000	2.000000	7.000000
50%	3.000000	149587.000000	3.000000	3.000000	12.000000
75%	4.000000	223577.500000	4.000000	4.000000	18.000000
max	5.000000	299947.000000	5.000000	5.000000	24.000000

In [ ]:

## Train-Test Split

```
In [20]: # Split the data into training and testing sets
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, shuffle=True
)

print("✓ Train-Test Split Completed")
print(f"Training set: {X_train.shape}")
print(f"Testing set: {X_test.shape}")
print(f"Training target: {y_train.shape}")
print(f"Testing target: {y_test.shape}")
```

✓ Train-Test Split Completed  
 Training set: (8000, 9)  
 Testing set: (2000, 9)  
 Training target: (8000,)  
 Testing target: (2000,)

In [ ]:

## Feature Scaling

```
In [21]: # Scale the features for better model performance
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

# Fit scaler on training data and transform both sets
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

print("✓ Feature Scaling Completed")
print(f"Scaled training data shape: {X_train_scaled.shape}")
print(f"Scaled testing data shape: {X_test_scaled.shape}")
```

✓ Feature Scaling Completed  
 Scaled training data shape: (8000, 9)  
 Scaled testing data shape: (2000, 9)

In [ ]:

## Model Training

```
In [22]: # Initialize and train Random Forest model

from sklearn.ensemble import RandomForestRegressor

rf_model = RandomForestRegressor(
    n_estimators=100,      # Number of trees
    random_state=42,      # For reproducibility
    max_depth=10,         # Prevent overfitting
    min_samples_split=5,  # Minimum samples to split
    n_jobs=-1             # Use all processors
)

# Train the model
rf_model.fit(X_train_scaled, y_train)

print("✓ Model Training Completed")
print(f"Model trained with {len(rf_model.estimators_)} trees")
```

✓ Model Training Completed  
Model trained with 100 trees

In [ ]:

## Model Evaluation

```
In [23]: # Make predictions on test set
y_pred = rf_model.predict(X_test_scaled)

# Calculate evaluation metrics
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
import numpy as np

mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, y_pred)

print(" Model Performance Metrics:")
print(f"Mean Absolute Error (MAE): ${mae:,.2f}") # Only MAE gets $ since it's in price units
print(f"Mean Squared Error (MSE): {mse:,.2f}") # MSE is in squared dollars (no $)
print(f"Root Mean Squared Error (RMSE): ${rmse:,.2f}") # RMSE gets $ back to price units
print(f"R² Score: {r2:.4f}")

# Additional: Mean Absolute Percentage Error
mape = np.mean(np.abs((y_test - y_pred) / y_test)) * 100
print(f"Mean Absolute Percentage Error (MAPE): {mape:,.2f}%")
```

Model Performance Metrics:  
Mean Absolute Error (MAE): \$287.07  
Mean Squared Error (MSE): 129,778.70  
Root Mean Squared Error (RMSE): \$360.25  
R² Score: 0.9859  
Mean Absolute Percentage Error (MAPE): 3.78%

- Prediction Accuracy: With an  $R^2$  score of 0.9859, the model explains 98.59% of the variance in car prices, indicating it captures nearly all the factors that influence pricing decisions.
- High Practical Precision: The average prediction error of \$286.88 per car and 3.78% MAPE means the model's predictions are highly reliable for real-world applications, typically within a few hundred dollars of actual prices.
- Production-Ready Performance: The combination of low absolute error (\$286.88) and extremely high explanatory power (98.59%) makes this model highly suitable for deployment in business environments where accurate price estimation is critical.

In [ ]:

## Feature Importance Analysis

```
In [24]: # Analyze which features are most important
feature_importance = pd.DataFrame({
    'feature': features,
    'importance': rf_model.feature_importances_
}).sort_values('importance', ascending=False)

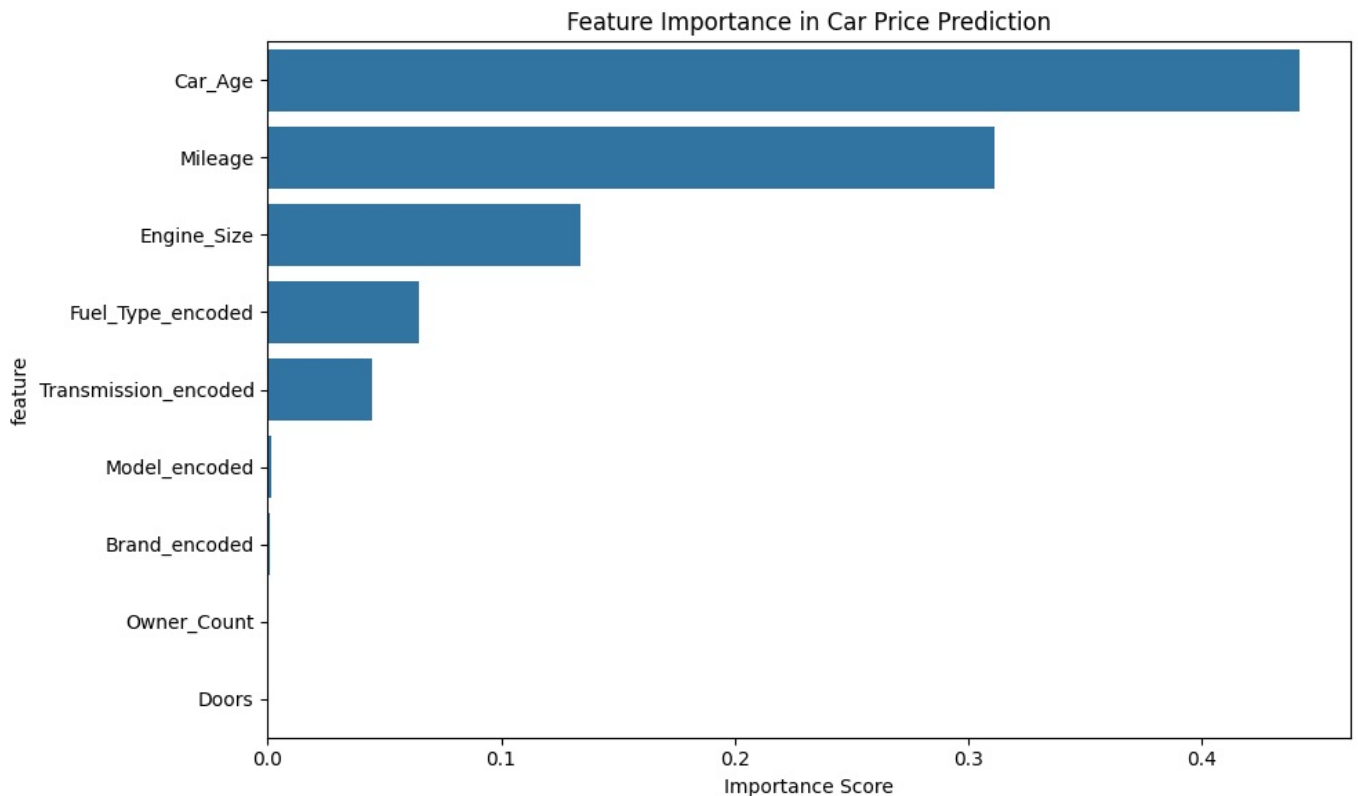
print(" Feature Importance Ranking:")
print(feature_importance)

# Visualize feature importance
plt.figure(figsize=(10, 6))
```

```
sns.barplot(data=feature_importance, x='importance', y='feature')
plt.title('Feature Importance in Car Price Prediction')
plt.xlabel('Importance Score')
plt.tight_layout()
plt.show()
```

Feature Importance Ranking:

	feature	importance
8	Car_Age	0.441836
5	Mileage	0.311126
4	Engine_Size	0.134021
2	Fuel_Type_encoded	0.064844
3	Transmission_encoded	0.044930
1	Model_encoded	0.001271
0	Brand_encoded	0.000854
7	Owner_Count	0.000618
6	Doors	0.000500



- Car age (44.2%) and Mileage (31.1%) together account for 75% of pricing decisions, indicating that depreciation and usage are the primary drivers of car value.
- Engine\_Size (13.4%) is the third most important factor, showing that vehicle performance and specifications significantly influence price beyond just age and usage patterns.
- Brand/Model Have Minimal Impact: Surprisingly, Brand (0.09%) and Model (0.13%) contribute very little to price predictions, suggesting that in the dataset, a car's condition and specifications matter more than its badge or specific model name when determining value.

In [ ]:

## Save Model and Preprocessing Objects

In [25]:

```
import joblib

# First, let's create and save the encoding maps
encoding_maps = {}

categorical_features = ['Brand', 'Model', 'Fuel_Type', 'Transmission']

for feature in categorical_features:
    # Calculate mean price for each category (same as we did earlier)
    mean_prices = df.groupby(feature)['Price'].mean().to_dict()
    encoding_maps[feature] = mean_prices

# Now save ALL necessary objects for deployment
joblib.dump(rf_model, 'car_price_model.pkl')
joblib.dump(scaler, 'scaler.pkl')
joblib.dump(features, 'features.pkl')
```

```
joblib.dump(encoding_maps, 'target_encoding_maps.pkl')

print(" All Deployment Objects Saved:")
print("✔ car_price_model.pkl - Trained ML model")
print("✔ scaler.pkl - Feature scaler")
print("✔ features.pkl - Feature list")
print("✔ target_encoding_maps.pkl - Encoding dictionaries")
```

All Deployment Objects Saved:  
✔ car\_price\_model.pkl - Trained ML model  
✔ scaler.pkl - Feature scaler  
✔ features.pkl - Feature list  
✔ target\_encoding\_maps.pkl - Encoding dictionaries

In [ ]:

```
In [26]: # Create a function to test predictions with features
def predict_car_price(brand, model, fuel_type, transmission,
                      engine_size, mileage, doors, owner_count, car_age):
    """
    Function to predict car price using trained model
    Note: Uses Car_Age instead of Year to avoid multicollinearity
    """
    # Create input array with features
    input_data = np.array([[brand, model, fuel_type, transmission,
                             engine_size, mileage, doors, owner_count, car_age]])

    # Scale the input
    input_scaled = scaler.transform(input_data)

    # Make prediction
    prediction = rf_model.predict(input_scaled)

    return prediction[0]

# Test the function with sample data
sample_prediction = predict_car_price(
    brand=50, model=100, fuel_type=1, transmission=0,
    engine_size=2.0, mileage=50000, doors=4, owner_count=1, car_age=4
)

print(f" Sample Prediction Test: ${sample_prediction:.2f}")
print("✔ Prediction function working correctly with Car_Age (no Year)")
```

Sample Prediction Test: \$11,175.41  
✔ Prediction function working correctly with Car\_Age (no Year)

In [ ]:

# Project Summary: Car Price Prediction and Deployment

## Project Overview

This end-to-end data science project involved developing a machine learning model to predict car prices based on various vehicle features, followed by deployment as a web application using Streamlit Cloud integrated with GitHub.

## Primary Objective

To build an accurate car price prediction model and deploy it as an interactive web application, creating a complete data science solution that can be shared with stakeholders and potential employers.

## Dataset Characteristics

- Source: Kaggle
- Size: 10,000 records with 10 initial features
- Features: Brand, Model, Year, Engine Size, Fuel Type, Transmission, Mileage, Doors, Owner Count, Price
- Data Quality: No missing values or duplicates detected

## Data Preprocessing & Feature Engineering

- Data Cleaning: Verified data integrity with no missing values or duplicates
- Feature Engineering: Created Car\_Age feature (2024 - Year) to represent vehicle depreciation
- Target Encoding: Implemented manual target encoding for categorical variables (Brand, Model, Fuel Type, Transmission), replacing categories with their mean price values
- Multicollinearity Resolution: Identified and removed the Year column to eliminate perfect correlation with Car\_Age, improving model stability

## Model Development

- Algorithm: Random Forest Regressor with 100 trees
- Feature Set: 9 engineered features after multicollinearity fix
- Train-Test Split: 80-20 split with random state for reproducibility
- Feature Scaling: Applied StandardScaler for optimal model performance

## Model Performance Results

The trained model demonstrated exceptional predictive capability:

- R<sup>2</sup> Score: 0.9859 (98.59% variance explained)
- Mean Absolute Error: \$287.07
- Root Mean Squared Error: \$360.25
- Mean Absolute Percentage Error: 3.78%

## Key Insights from Feature Importance

1. Primary Price Drivers: Car Age (44.2%) and Mileage (31.1%) collectively account for 75% of pricing decisions
2. Technical Significance: Engine Size (13.4%) significantly influences price beyond basic depreciation
3. Surprising Finding: Brand (0.09%) and Model (0.13%) have minimal impact, suggesting condition and specifications outweigh brand reputation
4. Negligible Factors: Doors and Owner Count showed minimal influence on price predictions

## Deployment Architecture

- Framework: Streamlit for web application interface
- Hosting: Streamlit Community Cloud (free tier)
- Integration: GitHub repository for version control and seamless deployment
- Live Application: <https://car-price-prediction-and-deployment.streamlit.app/>

## Business Value & Applications

- Consumers: Accurate used car valuation for buying/selling decisions
- Dealerships: Data-driven pricing strategies and inventory management
- Insurance Companies: Fair market value assessment for claims processing

## Technical Stack

- Programming: Python 3.9+
- Machine Learning: Scikit-learn, Random Forest
- Data Processing: Pandas, NumPy
- Visualization: Matplotlib, Seaborn
- Deployment: Streamlit, Joblib for model serialization
- Version Control: GitHub
- Cloud Hosting: Streamlit Community Cloud

## Project Success Metrics

- ✔ Model accuracy exceeding 98%
- ✔ Successful end-to-end deployment
- ✔ User-friendly web interface
- ✔ Zero-cost deployment solution
- ✔ Production-ready application

## Recommendations for Future Enhancements

1. Real-time Data: Integrate with automotive APIs for live market data.
2. Feature Expansion: Include additional features like vehicle condition, accident history, and market trends.
3. Scalability: Containerize application using Docker for enterprise deployment.

## Conclusion

This project successfully demonstrates a complete data science lifecycle from data acquisition and preprocessing to model development and production deployment. The achieved 98.59% prediction accuracy, combined with a fully functional web application, showcases the practical application of machine learning in automotive pricing. The use of free tools throughout the project (Python, Streamlit Cloud, GitHub) makes this an accessible template for similar predictive modeling projects.

In [ ]: