

This process will include data loading, cleaning, visualization, and model creation.

1. Define the Problem and Objective

- **Objective:** Predict the price of a car based on various features such as brand, model, year of manufacture, mileage, engine size, etc.
 - **Goal:** Perform EDA to understand the relationships between features and the target variable (`price`), which represents the price of the car.
-

2. Load the Dataset

```
import pandas as pd

# Load Car Price dataset
df = pd.read_csv('car_price.csv')

# Display the first few rows of the dataset
df.head()
```

3. Data Exploration (EDA)

Check the Shape and Data Types

```
# Check dataset shape
print(f"Dataset Shape: {df.shape}")

# Data types and missing values
print(df.info())

# Summary statistics for numerical columns
print(df.describe())
```

Visualize Distributions

```
import matplotlib.pyplot as plt
import seaborn as sns

# Distribution of Car Price
sns.histplot(df['price'], kde=True)
plt.title("Car Price Distribution")
plt.show()

# Distribution of Year of Manufacture
sns.histplot(df['year'], kde=True)
plt.title("Car Manufacture Year Distribution")
plt.show()
```

4. Data Cleaning

Handle Missing Values

```
# Check for missing values
print(df.isnull().sum())

# Handle missing values (if any)
# For example, fill missing values in 'mileage' with the median value
df['mileage'].fillna(df['mileage'].median(), inplace=True)

# Check if missing values are handled
print(df.isnull().sum())
```

Handle Duplicates

```
# Remove duplicate rows if any
df.drop_duplicates(inplace=True)
```

5. Feature Engineering (Optional)

Create New Features (if needed)

```
# Example: Create a new feature 'age_of_car' by calculating the difference
between the current year and the car's manufacturing year
df['age_of_car'] = 2025 - df['year']

# Display first few rows to check new feature
print(df[['year', 'age_of_car']].head())
```

6. Data Preprocessing

Convert Categorical Variables to Numerical

```
# Convert 'brand' or 'model' (if categorical) into numerical using one-hot
encoding
df = pd.get_dummies(df, columns=['brand', 'model'], drop_first=True)

# Example: Convert 'fuel_type' to numerical values
df['fuel_type'] = df['fuel_type'].map({'Petrol': 0, 'Diesel': 1, 'Electric':
2})

# Check the changes in data
print(df.head())
```

Scale Numerical Features

```
from sklearn.preprocessing import StandardScaler

# Initialize the scaler
scaler = StandardScaler()

# Scale the numerical features (e.g., price, mileage, etc.)
scaled_features = df[['mileage', 'price', 'engine_size', 'age_of_car']] #
Specify numerical columns to scale
df[['mileage', 'price', 'engine_size', 'age_of_car']] =
scaler.fit_transform(scaled_features)
```

Split the Data into Features (X) and Target (y)

```
# Define features (X) and target (y)
X = df.drop(columns=['price']) # Exclude target variable (price)
y = df['price']
```

7. Split the Data into Training and Testing Sets

```
from sklearn.model_selection import train_test_split

# Split data into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

8. Model Selection

Choose the Model

```
from sklearn.ensemble import RandomForestRegressor

# Initialize the Random Forest model
model = RandomForestRegressor(n_estimators=100, random_state=42)
```

9. Model Training

```
# Train the model on the training data
model.fit(X_train, y_train)
```

10. Model Evaluation

Predictions and Evaluation Metrics

```
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

# Make predictions
y_pred = model.predict(X_test)
```

```
# Evaluate the model
mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f'Mean Absolute Error: {mae}')
print(f'Mean Squared Error: {mse}')
print(f'R-squared: {r2}')
```

Residuals Plot

```
# Plot residuals to understand prediction errors
sns.scatterplot(x=y_test, y=y_test - y_pred)
plt.title("Residuals Plot")
plt.xlabel("True Values")
plt.ylabel("Residuals")
plt.show()
```

11. Model Interpretation and Insights

Feature Importance

```
# Get feature importances from the trained model
importances = model.feature_importances_

# Create a DataFrame to display features and their importance
features = X.columns
feature_importance_df = pd.DataFrame({'Feature': features, 'Importance':
importances})

# Sort by importance
feature_importance_df = feature_importance_df.sort_values(by='Importance',
ascending=False)

# Plot feature importances
sns.barplot(x='Importance', y='Feature', data=feature_importance_df)
plt.title("Feature Importance")
plt.show()

print(feature_importance_df)
```

12. Fine-Tuning and Optimization (Optional)

Hyperparameter Tuning

```
from sklearn.model_selection import GridSearchCV

# Set up grid search parameters for Random Forest Regressor
param_grid = {
    'n_estimators': [100, 200, 300],
```

```
'max_depth': [None, 10, 20],
'min_samples_split': [2, 5, 10]
}

# Initialize GridSearchCV
grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=3,
n_jobs=-1, verbose=2)

# Fit grid search
grid_search.fit(X_train, y_train)

# Get best parameters
print("Best Parameters:", grid_search.best_params_)
```

13. Final Reporting

Summarize Insights

- **EDA Insights:** Discuss the correlation between car features like mileage, engine size, and the car price. Highlight any trends in car prices and any relationships found between features.
 - **Model Evaluation:** Summarize the Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared values.
 - **Important Features:** Identify the most important features influencing car price predictions (e.g., age, mileage, engine size).
 - **Next Steps:** Explore other models (e.g., Gradient Boosting) or perform hyperparameter tuning for better results.
-

Additional Notes:

- Start with a **baseline model** like Random Forest Regressor and evaluate its performance.
 - Use **cross-validation** to assess model robustness.
 - Explore **feature engineering** like transforming variables or combining features for better prediction accuracy.
-