

```
In [12]: pip install numpy pandas matplotlib seaborn plotly

Requirement already satisfied: numpy in c:\users\laxay\python\python3\lib\site-packages (1.26.4)
Requirement already satisfied: pandas in c:\users\laxay\python\python3\lib\site-packages (2.2.2)
Requirement already satisfied: matplotlib in c:\users\laxay\python\python3\lib\site-packages (3.8.0)
Requirement already satisfied: seaborn in c:\users\laxay\python\python3\lib\site-packages (0.13.2)
Requirement already satisfied: plotly in c:\users\laxay\python\python3\lib\site-packages (5.22.0)
Requirement already satisfied: python-dateutil<2.8.2 in c:\users\laxay\python\python3\lib\site-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz<2023.1 in c:\users\laxay\python\python3\lib\site-packages (from pandas) (2024.2)
Requirement already satisfied: contourpy<1.1 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (1.1.1)
Requirement already satisfied: cycler<0.12 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools<4.5 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (4.58.0)
Requirement already satisfied: kiwisolver<3.1 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (3.1.0)
Requirement already satisfied: pillow<10 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (10.4.0)
Requirement already satisfied: pyparsing<3.1 in c:\users\laxay\python\python3\lib\site-packages (from matplotlib) (3.1.2)
Requirement already satisfied: fontcabin in c:\users\laxay\python\python3\lib\site-packages (from plotly) (1.6.0)
Requirement already satisfied: fontcabin in c:\users\laxay\python\python3\lib\site-packages (from plotly) (1.6.0)
Note: you may need to restart the kernel to use updated packages.

In [13]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
import warnings
warnings.filterwarnings("ignore")

In [15]: df = pd.read_csv('car_data.csv')

In [16]: print("First few rows of the dataset:")
df.head()

Out[16]:
   Car_Name  Year  Selling_Price  Present_Price  Driven_kms  Fuel_Type  Selling_Type  Transmission  Owner
0      rzt  2014         3.35          5.59      27000      Petrol      Dealer      Manual      0
1      sz4  2013         4.75          9.54      43000      Diesel      Dealer      Manual      0
2      oaz  2017         7.25          9.85      6900      Petrol      Dealer      Manual      0
3  wagon r  2011         2.85          4.15      5200      Petrol      Dealer      Manual      0
4      owt  2014         4.60          6.87      42400      Diesel      Dealer      Manual      0

In [17]: df.columns

Out[17]:
Index(['Car_Name', 'Year', 'Selling_Price', 'Present_Price', 'Driven_kms',
      'Fuel_Type', 'Selling_Type', 'Transmission', 'Owner'],
      dtype='object')

In [18]: df.shape

Out[18]:
(361, 9)

In [19]: df.describe().style.format(precision=2).background_gradient(cmap='Blues')

Out[19]:
   Car_Name  Year  Selling_Price  Present_Price  Driven_kms  Owner
count  361.00         301.00         501.00         501.00         301.00
mean    2014.63          4.66          7.63      35547.21          0.04
std       2.89          5.08          8.64      36688.85          0.25
min    2003.00          0.10          0.32          600.00          0.00
25%  2012.00          0.90          1.20      15000.00          0.00
50%  2014.00          3.60          6.40      32000.00          0.00
75%  2016.00          6.00          9.90      48767.00          0.00
max    2016.00          6.00          9.90      50000.00          0.00

In [20]: df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 361 entries, 0 to 360
Data columns (total 9 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Car_Name    361 non-null     object
 1   Year        361 non-null     int64
 2   Selling_Price  361 non-null     float64
 3   Present_Price  361 non-null     float64
 4   Driven_kms   361 non-null     int64
 5   Fuel_Type    361 non-null     object
 6   Selling_Type  361 non-null     object
 7   Transmission  361 non-null     object
 8   Owner       361 non-null     int64
dtypes: float64(2), int64(3), object(4)
memory usage: 21.2+ KB

In [21]: # Check Missing Values
print("Missing Values:")
df.isnull().sum()

Missing Values:
Car_Name      0
Year          0
Selling_Price  0
Present_Price  0
Driven_kms     0
Fuel_Type     0
Selling_Type  0
Transmission   0
Owner         0
dtype: int64

In [22]: # Data types of columns
print("Data types of columns:")
print(df.dtypes)

Data types of columns:
Car_Name      object
Year          int64
Selling_Price float64
Present_Price float64
Driven_kms    int64
Fuel_Type     object
Selling_Type  object
Transmission  object
Owner         int64

In [23]: # Check for duplicate values
print("Duplicate Values:")
df.duplicated().sum()

Duplicate values:
2

In [24]: # Drop duplicate values
df = df.drop_duplicates()
df.duplicated().sum()

0

In [25]: # Check for outliers using boxplots
plt.figure(figsize=(10, 6))
sns.boxplot(df['Selling_Price'], palette='pastel')
plt.title('Boxplot of Selling Price')
plt.show()

Boxplot of Selling Price

In [26]: # Select numerical columns
numerical_columns = ['Year', 'Selling_Price', 'Present_Price', 'Driven_kms', 'Owner']

# Create a DataFrame containing only the numerical columns
numerical_df = df[numerical_columns]

# Calculate the correlation matrix for numerical columns
correlation_matrix = numerical_df.corr()

# Create a heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(numerical_df, annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Heatmap (Numerical Columns)')
plt.show()

Correlation Heatmap (Numerical Columns)

In [27]: numerical_features = ['Year', 'Driven_kms', 'Selling_Price', 'Present_Price']

# Feature to numerical features
plt.figure(figsize=(10, 6))
sns.distplot(numerical_df[numerical_features])
plt.title('Distribution of Features')
plt.show()

Distribution of Year

In [28]: # Figure size 1000x600 with 8 Axes>

Count
0
10
20
30
40
50
60
2004 2006 2008 2010 2012 2014 2016 2018

Distribution of Driven_kms

In [29]: # Figure size 1000x600 with 8 Axes>

Count
0
10
20
30
40
50
60
0 100000 200000 300000 400000 500000

Distribution of Selling_Price

In [30]: # Figure size 1000x600 with 8 Axes>

Count
0
20
40
60
80
100
0 20 40 60 80
Present_Price

Distribution of Present_Price

In [31]: # Scatter plots
plt.figure(figsize=(12, 8))
sns.scatterplot(x=Driven_kms, y=Selling_Price, data=df)
plt.title('Relationship between Driven_kms and Selling_Price')
plt.show()

Relationship between Driven_kms and Selling_Price

In [32]: plt.figure(figsize=(12, 8))
sns.scatterplot(x=Present_Price, y=Selling_Price, data=df)
plt.title('Relationship between Present_Price and Selling_Price')
plt.xlabel('Present_Price')
plt.ylabel('Selling_Price')
plt.show()

Relationship between Present_Price and Selling_Price

In [33]: plt.figure(figsize=(12, 8))
sns.scatterplot(x=Present_Price, y=Selling_Price, data=df)
plt.title('Relationship between Present_Price and Selling_Price')
plt.xlabel('Present_Price')
plt.ylabel('Selling_Price')
plt.show()

Relationship between Present_Price and Selling_Price

In [34]: # Define a color palette for the plots
sns.set_palette('magma')

categorical_features = ['Fuel_Type', 'Selling_Type', 'Transmission', 'Owner']

# For each categorical feature:
plt.figure(figsize=(10, 6))
sns.countplot(x=feature, data=df, palette=palette)
plt.title(f'Frequency of {feature}')
plt.xlabel(f'{feature}')
plt.ylabel(f'Frequency')
plt.grid(True)
plt.show()

Frequency of Fuel_Type

In [35]: # Frequency of Fuel_Type
plt.figure(figsize=(10, 6))
sns.countplot(x='Fuel_Type', data=df)
plt.title('Frequency of Fuel_Type')
plt.xlabel('Fuel_Type')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Fuel_Type

In [36]: # Frequency of Selling_Type
plt.figure(figsize=(10, 6))
sns.countplot(x='Selling_Type', data=df)
plt.title('Frequency of Selling_Type')
plt.xlabel('Selling_Type')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Selling_Type

In [37]: # Frequency of Transmission
plt.figure(figsize=(10, 6))
sns.countplot(x='Transmission', data=df)
plt.title('Frequency of Transmission')
plt.xlabel('Transmission')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Transmission

In [38]: # Frequency of Owner
plt.figure(figsize=(10, 6))
sns.countplot(x='Owner', data=df)
plt.title('Frequency of Owner')
plt.xlabel('Owner')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Owner

In [39]: df.columns

Out[39]:
Index(['Car_Name', 'Year', 'Selling_Price', 'Present_Price', 'Driven_kms',
      'Fuel_Type', 'Selling_Type', 'Transmission', 'Owner'],
      dtype='object')

In [40]: n = 20 # Number of top car models to plot
top_car_models = df[['Car_Name', 'Selling_Price']].nlargest(n)

plt.figure(figsize=(10, 6))
sns.barplot(x=top_car_models.index, y=top_car_models.Selling_Price)
plt.title('Top 20 Car Models by Frequency')
plt.xlabel('Car Model')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Top 20 Car Models by Frequency

In [41]: # Calculate average price for each car model
avg_prices_by_car = df.groupby('Car_Name')['Selling_Price'].mean().sort_values(ascending=False)

# Plot top 8 car models by average price
n = 8 # Number of top car models to plot
top_car_models = avg_prices_by_car.head(n)

plt.figure(figsize=(10, 6))
sns.barplot(x=top_car_models.index, y=top_car_models.Selling_Price)
plt.title('Top 8 Car Models by Average Price')
plt.xlabel('Car Model')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Top 8 Car Models by Average Price

In [42]: # Frequency of Fuel_Type
plt.figure(figsize=(10, 6))
sns.countplot(x='Fuel_Type', data=df)
plt.title('Frequency of Fuel_Type')
plt.xlabel('Fuel_Type')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Fuel_Type

In [43]: # Frequency of Selling_Type
plt.figure(figsize=(10, 6))
sns.countplot(x='Selling_Type', data=df)
plt.title('Frequency of Selling_Type')
plt.xlabel('Selling_Type')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Selling_Type

In [44]: # Frequency of Transmission
plt.figure(figsize=(10, 6))
sns.countplot(x='Transmission', data=df)
plt.title('Frequency of Transmission')
plt.xlabel('Transmission')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Transmission

In [45]: # Frequency of Owner
plt.figure(figsize=(10, 6))
sns.countplot(x='Owner', data=df)
plt.title('Frequency of Owner')
plt.xlabel('Owner')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()

Frequency of Owner

In [46]: # Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

In [47]: # Fit the Linear Regression model
model = LinearRegression()
model.fit(X_train, y_train)

In [48]: # Evaluate the model
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error: ", mse)

Mean Squared Error: 6.575410693250
R-squared: 0.74993435643088

In [49]: # Make predictions
new_car = {'Year': 2000, 'Driven_kms': 0, 'Fuel_Type': 'Petrol', 'Selling_Type': 'Manual', 'Transmission': 'Manual'}
predicted_price = model.predict(new_car)
print("Predicted Selling Price: ", predicted_price[0])
```


