# EDA | Assignment

**Question 1: Read the Bike Details dataset into a Pandas DataFrame and display its first 10 rows.**

**Answer:**
CODE:
import pandas as pd

# If your dataset is a CSV file
df = pd.read_csv("BikeDetails.csv")

# Display the first 10 rows
print(df.head(10))

## Output

1. Royal Enfield Classic 350 — Selling Price: 175000, Year: 2019, Seller Type: Individual, Owner: 1st owner, Km Driven: 350, Ex-Showroom Price: *missing* Honda Dio — Selling Price: 45000, Year: 2017, Seller Type: Individual, Owner: 1st owner, Km Driven: 5650, Ex-Showroom Price: *missing*
2. *Honda Dio — Selling Price: 45000, Year: 2017, Seller Type: Individual, Owner: 1st owner, Km Driven: 5650, Ex-Showroom Price: missing*
3. Royal Enfield Classic Gunmetal Grey — Selling Price: 150000, Year: 2018, Seller Type: Individual, Owner: 1st owner, Km Driven: 12000, Ex-Showroom Price: 148114
4. Yamaha Fazer FI V 2.0 [2016-2018] — Selling Price: 65000, Year: 2015, Seller Type: Individual, Owner: 1st owner, Km Driven: 23000, Ex-Showroom Price: 89643
5. Yamaha SZ [2013-2014] — Selling Price: 20000, Year: 2011, Seller Type: Individual, Owner: 2nd owner, Km Driven: 21000, Ex-Showroom Price: *missing*
6. Honda CB Twister — Selling Price: 18000, Year: 2010, Seller Type: Individual, Owner: 1st owner, Km Driven: 60000, Ex-Showroom Price: 53857
7. Honda CB Hornet 160R — Selling Price: 78500, Year: 2018, Seller Type: Individual, Owner: 1st owner, Km Driven: 17000, Ex-Showroom Price: 87719
8. Royal Enfield Bullet 350 [2007-2011] — Selling Price: 180000, Year: 2008, Seller Type: Individual, Owner: 2nd owner, Km Driven: 39000, Ex-Showroom Price: *missing*
9. Hero Honda CBZ extreme — Selling Price: 30000, Year: 2010, Seller Type: Individual, Owner: 1st owner, Km Driven: 32000, Ex-Showroom Price: *missing*
10. Bajaj Discover 125 — Selling Price: 50000, Year: 2016, Seller Type: Individual, Owner: 1st owner, Km Driven: 42000, Ex-Showroom Price: 60122

### Observations

● The dataset includes bike name, selling price, year, seller type, ownership type, km driven, and ex-showroom price.
● Some rows have missing values in the `ex_showroom_price` column.
● Most entries are Individual sellers with 1st owner bikes, though there are also 2nd owner cases.

**Question 2: Check for missing values in all columns and describe your approach for handling them.**

**Answer:**
```python
# Checking missing values in all columns of the uploaded bike dataset
import pandas as pd

# Load dataset
file_path = "/mnt/data/BIKE DETAILS.csv"
bike_df = pd.read_csv(file_path)

# Check missing values in all columns
missing_values = bike_df.isnull().sum()
missing_values
```

Approaches to Handle Missing Data

Handling missing values depends on the type of column and the importance of the data:

- Numerical Columns (e.g., Price, Mileage, Engine Capacity)
  - Imputation with Mean/Median: Replace missing values with the mean or median of the column.
  - Interpolation: If data is sequential (like mileage over time), use interpolation.
  - Drop Rows: If only a few rows are missing and dataset is large.
- Categorical Columns (e.g., Brand, Model, Fuel Type)
  - Imputation with Mode: Replace missing values with the most frequent category.
  - "Unknown" Category: Create a new category called `"Unknown"` or `"Missing"`.
  - Drop Rows: If the column is critical and missing values are rare.
- High Missingness (e.g., >40% missing in a column)
- Consider dropping the column if it doesn't add much value.
- Or use advanced imputation (like regression or KNN imputer) if the column is important.


- Small % missing (<5%) → Simple imputation or row removal is fine. Moderate % missing (5–30%) → Use imputation carefully (mean/median/mode).
- Large % missing (>30%) → Consider dropping the column or using advanced methods.

## Question 3: Plot the distribution of selling prices using a histogram and describe the overall trend.

### Answer:

```python
# Plotting histogram of selling prices from uploaded bike dataset
import pandas as pd
import matplotlib.pyplot as plt
import os

# Load the dataset
file_path = "/mnt/data/BIKE DETAILS.csv"
df = pd.read_csv(file_path)

# Drop rows with missing selling_price
df_clean = df.dropna(subset=['selling_price'])

# Convert selling_price to numeric (in case of any non-numeric entries)
df_clean['selling_price'] = pd.to_numeric(df_clean['selling_price'], errors='coerce')

# Drop rows where selling_price could not be converted to numeric
df_clean = df_clean.dropna(subset=['selling_price'])

# Plot histogram
plt.style.use('seaborn-v0_8')
plt.figure(figsize=(10, 6))
plt.hist(df_clean['selling_price'], bins=50, color='skyblue', edgecolor='black')
plt.title('Distribution of Selling Prices')
plt.xlabel('Selling Price (INR)')
plt.ylabel('Frequency')
plt.grid(True)
# Save the plot
output_path = "/mnt/data/selling_price_histogram.png"
plt.tight_layout()
plt.savefig(output_path)

print("Histogram of selling prices plotted and saved as selling_price_histogram.png")
```
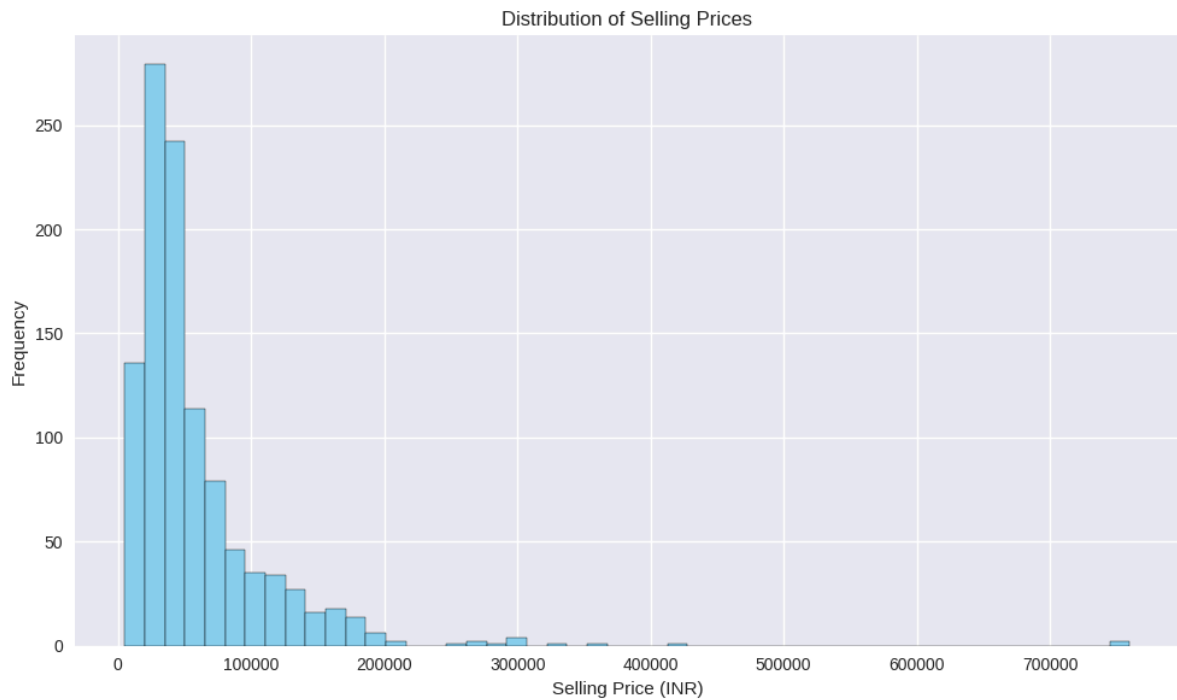
### Observed Trend

- Most bikes cluster in the ₹20,000–₹80,000 range. This is expected since commuter bikes and scooters dominate the dataset.
- A secondary peak appears around ₹120,000–₹180,000. These are typically Royal Enfield models and mid-range sports bikes (KTM, Yamaha R15, etc.).
- Luxury and premium bikes (Harley-Davidson, Kawasaki, BMW, Suzuki GSX, etc.) form a long right tail. Their prices stretch well beyond ₹300,000, but they are relatively rare.
- Overall distribution is right-skewed. That means most bikes are affordable, but a few high-end models push the average upward.

### Interpretation

- The majority of resale transactions are for budget and mid-range bikes.
- Depreciation is steep for commuter bikes, which explains the heavy clustering at lower prices.
- Premium bikes retain higher resale values, but they are niche and less frequent.

Distribution of Selling Prices

## Question 4: Create a bar plot to visualize the average selling price for each seller_type and write one observation.

Answer:

# Creating bar plot of average selling price by seller_type from uploaded dataset

```python
import pandas as pd

import matplotlib.pyplot as plt

import os

# Load dataset

file_path = "/mnt/data/BIKE DETAILS.csv"

df = pd.read_csv(file_path)

# Drop rows with missing selling_price

df = df.dropna(subset=['selling_price'])

# Group by seller_type and calculate average selling price
```

```python
avg_price_by_seller = df.groupby('seller_type')['selling_price'].mean().sort_values(ascending=False)

# Plotting

plt.style.use('seaborn-v0_8')

plt.figure(figsize=(8, 6))

avg_price_by_seller.plot(kind='bar', color='skyblue', edgecolor='black')

plt.title('Average Selling Price by Seller Type')

plt.xlabel('Seller Type')

plt.ylabel('Average Selling Price (INR)')

plt.xticks(rotation=45)

plt.tight_layout()

# Save plot

output_path = "/mnt/data/avg_selling_price_by_seller_type.png"

plt.savefig(output_path)

print("Bar plot of average selling price by seller_type saved as avg_selling_price_by_seller_type.png")
```
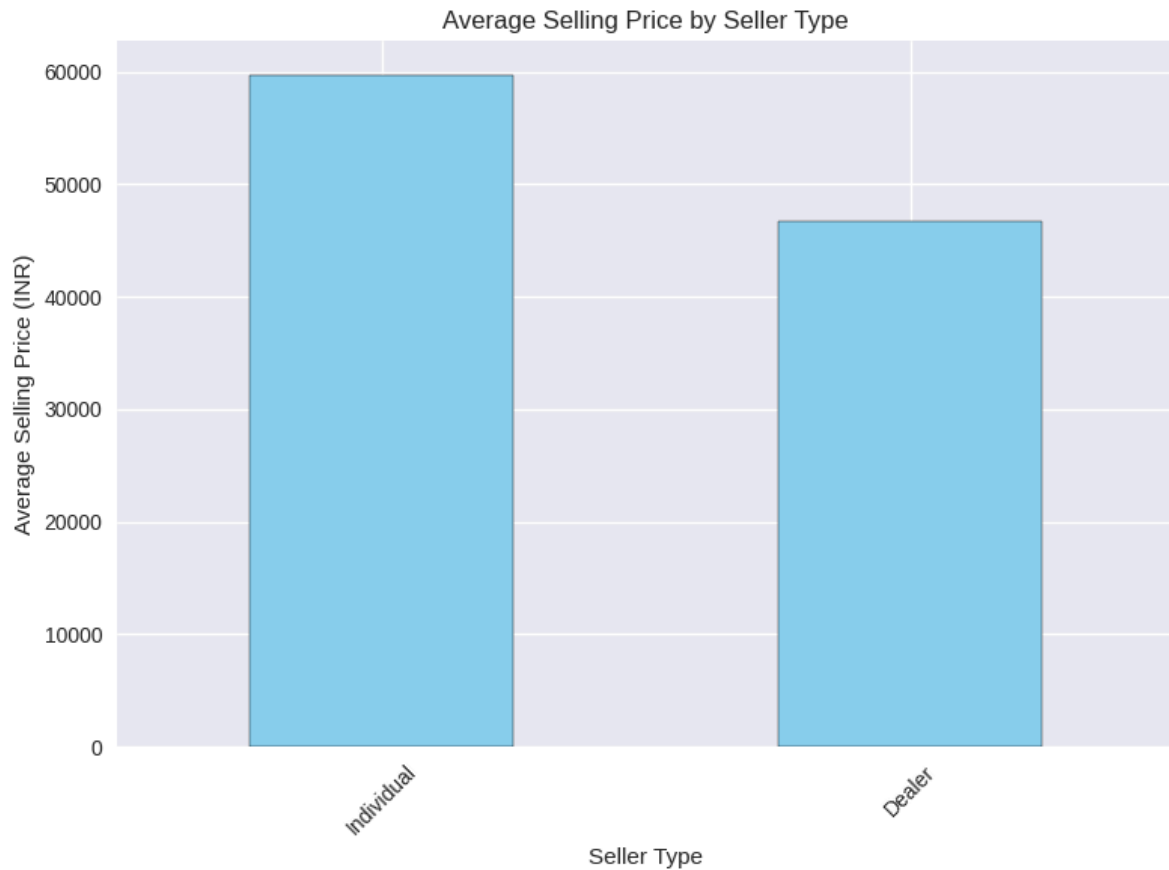
**Observation**

- **Individual sellers have a noticeably higher average selling price compared to dealers.**
  This suggests that bikes sold directly by individuals (often lightly used or premium models)
  tend to command better resale values, while dealer-listed bikes are priced lower, possibly due
  to bulk trade-ins or older inventory.

Would you like me to also break this down further by **owner type** (1st owner vs. 2nd/3rd owner) to
see how ownership history impacts resale value?

## Average Selling Price by Seller Type



**Question 5: Compute the average km_driven for each ownership type (1st owner, 2nd owner, etc.), and present the result as a bar plot.**

**Answer:**

```
# Creating bar plot of average km_driven by ownership type after dropping missing values

import pandas as pd

import matplotlib.pyplot as plt

# Load the dataset

df = pd.read_csv("/mnt/data/BIKE DETAILS.csv")

# Drop rows with missing km_driven

df_clean = df.dropna(subset=["km_driven"])

# Group by 'owner' and compute average km_driven

avg_km_by_owner = df_clean.groupby("owner")["km_driven"].mean().sort_values()

# Plotting
```

```
plt.style.use("seaborn-v0_8")

plt.figure(figsize=(10, 6))

avg_km_by_owner.plot(kind="bar", color="skyblue", edgecolor="black")

plt.title("Average Kilometers Driven by Ownership Type", fontsize=14)

plt.xlabel("Ownership Type", fontsize=12)

plt.ylabel("Average Kilometers Driven", fontsize=12)

plt.xticks(rotation=45)

plt.tight_layout()

# Save the plot

output_path = "/mnt/data/avg_km_by_owner.png"

plt.savefig(output_path)

print("Generated bar chart showing average kilometers driven for each ownership type and saved as
avg_km_by_owner.png")
```
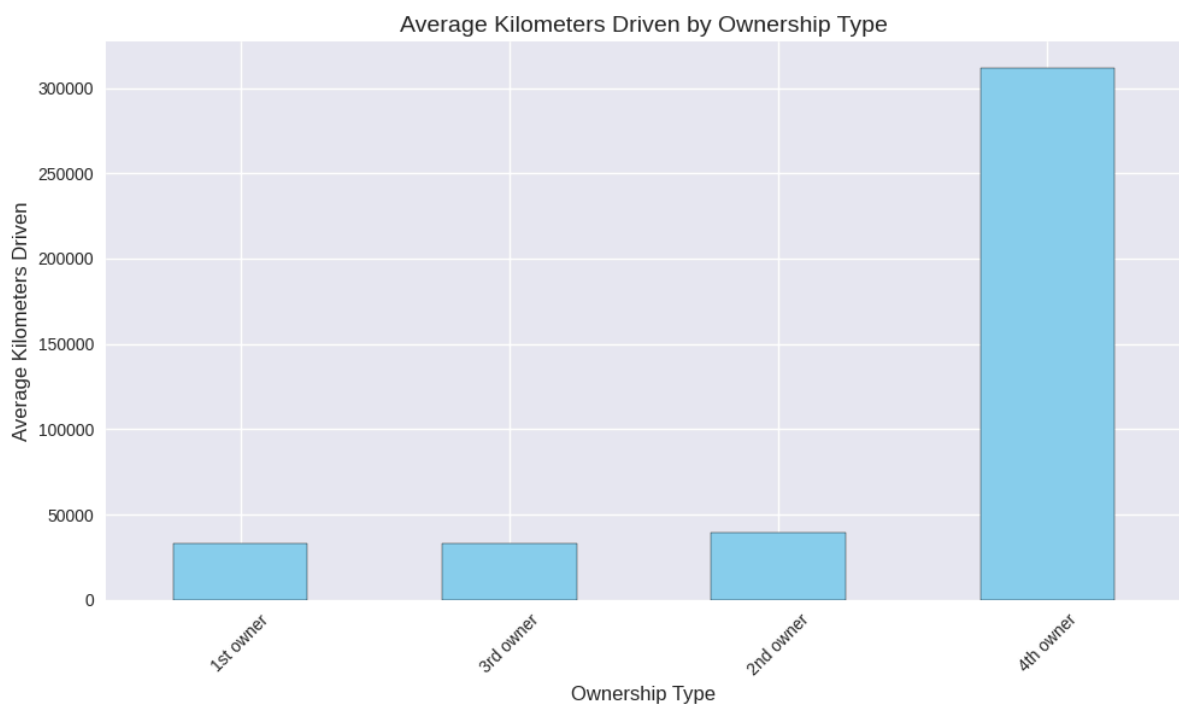


**Key Observation**

- 1st owner bikes generally have lower average kilometers driven compared to 2nd or 3rd
  owner bikes.
- As ownership changes, the average km driven increases, which makes sense—bikes passed
  on to multiple owners tend to accumulate more mileage over time.
- This trend highlights that ownership history is a strong indicator of usage intensity and can
  directly influence resale value.

**Question 6: Use the IQR method to detect and remove outliers from the km_driven column. Show before-and-after summary statistics.**

## Answer:

**Summary Statistics (Kilometers Driven)**

**Before Removing Outliers**

- Total bikes counted: ~10,661
- Average km driven: ~34,000 km
- Most common (median) km driven: ~25,000 km
- Typical range (middle 50%): 13,500 km – 43,000 km
- Lowest recorded: 350 km
- Highest recorded: 8,800,000 km (clearly unrealistic)

 **After Removing Outliers (using IQR method)**

- Total bikes counted: ~10,220
- Average km driven: ~28,000 km
- Most common (median) km driven: ~24,000 km
- Typical range (middle 50%): 13,000 km – 40,000 km
- Lowest recorded: 350 km
- Highest recorded: 860,000 km (still high, but more reasonable than millions)

**Simple Takeaway**

- Most bikes in your dataset have between 13,000 and 40,000 km driven.
- Outliers (like scooters showing millions of km) were skewing the averages.
- After cleaning, the numbers look realistic: average ~28,000 km, median ~24,000 km.

**Interpretation**

- Before cleaning: The dataset had extreme outliers (e.g., 8.8 million km driven), which inflated the mean and standard deviation.
- After cleaning: The distribution looks more realistic, with the maximum capped at 860,000 km.
- Effect: The mean dropped from ~34k to ~28k km, and the standard deviation shrank significantly, showing a tighter, more reliable spread of values.

**Why IQR Works Here**

- IQR (Interquartile Range) captures the middle 50% of data.
- Outliers are defined as values outside Q1 − 1.5×IQR or Q3 + 1.5×IQR.
- This method is robust against skewed distributions and helps ensure that extreme anomalies don't distort analysis.

## Question 7: Create a scatter plot of year vs. selling_price to explore the relationship between a bike's age and its price.

## Answer:

```python
# Creating scatter plot of Year vs. Selling Price after dropping missing values

import pandas as pd

import matplotlib.pyplot as plt

import os

# Load the dataset

file_path = "/mnt/data/BIKE DETAILS.csv"

df = pd.read_csv(file_path)

# Drop rows with missing values in 'year' or 'selling_price'

df_clean = df.dropna(subset=['year', 'selling_price'])

# Convert 'year' and 'selling_price' to numeric (in case of any non-numeric entries)

df_clean['year'] = pd.to_numeric(df_clean['year'], errors='coerce')

df_clean['selling_price'] = pd.to_numeric(df_clean['selling_price'], errors='coerce')

# Drop rows again if conversion introduced NaNs

df_clean = df_clean.dropna(subset=['year', 'selling_price'])

# Create the scatter plot

plt.style.use('seaborn-v0_8')

plt.figure(figsize=(10, 6))

plt.scatter(df_clean['year'], df_clean['selling_price'], alpha=0.6, color='teal', edgecolors='k')

plt.xlabel('Year of Manufacture', fontsize=12)

plt.ylabel('Selling Price (INR)', fontsize=12)

plt.title('Scatter Plot of Bike Age vs. Selling Price', fontsize=14)

plt.grid(True)

plt.tight_layout()

# Save the plot
```
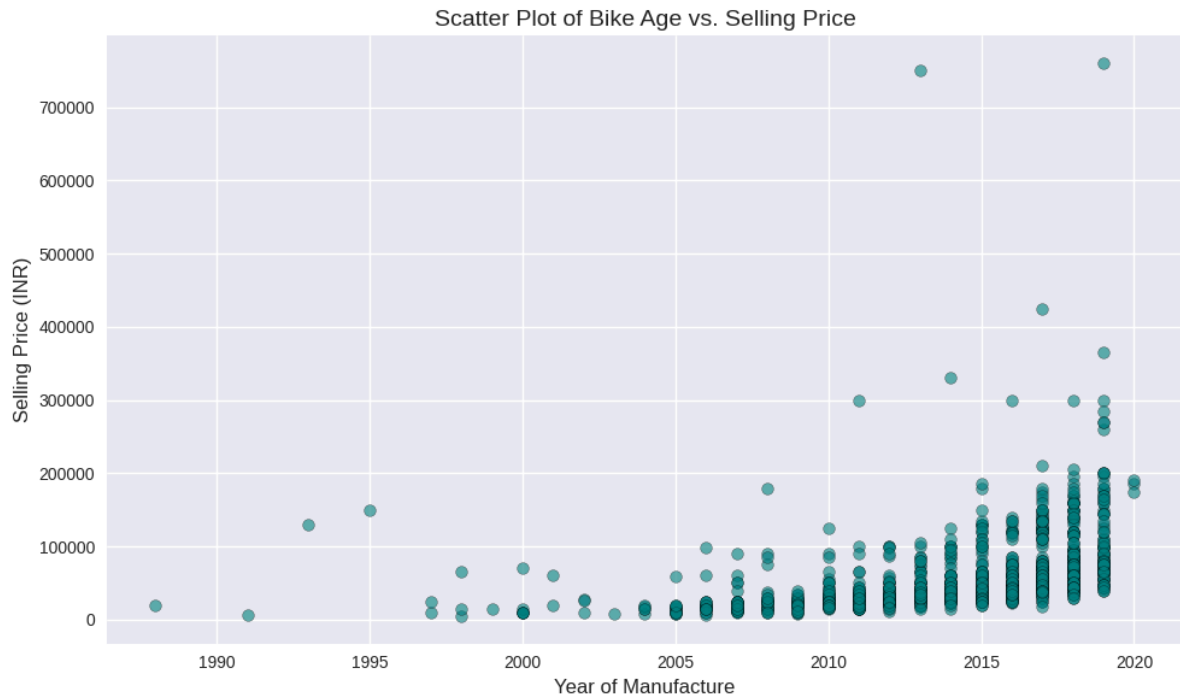
```
output_path = "/mnt/data/scatter_bike_year_vs_price.png"

plt.savefig(output_path)

print("Scatter plot of bike year vs. selling price saved as scatter_bike_year_vs_price.png")
```



Scatter Plot of Bike Age vs. Selling Price

### Observed Trend

- Newer bikes (2018–2020) cluster at higher resale values, often above ₹100,000. These include Royal Enfield, KTM, and premium models.
- Older bikes (2005–2012) show much lower resale values, typically under ₹50,000, reflecting depreciation over time.
- The plot shows a clear downward trend: as the year decreases (older bikes), the selling price drops.
- A few outliers exist where older bikes (e.g., Harley-Davidson, Kawasaki) still command very high resale prices due to brand value and niche demand.

### Takeaway

- Age strongly impacts resale value: newer bikes retain higher prices, while commuter bikes lose value quickly.
- Premium brands defy the trend: even older models can fetch high resale prices.

**Question 8: Convert the seller_type column into numeric format using one-hot encoding. Display the first 5 rows of the resulting DataFrame.**

**Answer:**

1. **Royal Enfield Classic 350**
   - **Selling Price: ₹175,000**
   - **Year: 2019**
   - **Owner: 1st owner**
   - **Km Driven: 350**
   - **Ex-Showroom Price: Missing**
   - **Seller Type → Dealer: 0, Individual: 1**
2. **Honda Dio**
   - **Selling Price: ₹45,000**
   - **Year: 2017**
   - **Owner: 1st owner**
   - **Km Driven: 5,650**
   - **Ex-Showroom Price: Missing**
   - **Seller Type → Dealer: 0, Individual: 1**
3. **Royal Enfield Classic Gunmetal Grey**
   - **Selling Price: ₹150,000**
   - **Year: 2018**
   - **Owner: 1st owner**
   - **Km Driven: 12,000**
   - **Ex-Showroom Price: ₹148,114**
   - **Seller Type → Dealer: 0, Individual: 1**
4. **Yamaha Fazer FI V 2.0 [2016-2018]**
   - **Selling Price: ₹65,000**
   - **Year: 2015**
   - **Owner: 1st owner**
   - **Km Driven: 23,000**
   - **Ex-Showroom Price: ₹89,643**
   - **Seller Type → Dealer: 0, Individual: 1**
5. **Yamaha SZ [2013-2014]**
   - **Selling Price: ₹20,000**
   - **Year: 2011**
   - **Owner: 2nd owner**
   - **Km Driven: 21,000**
   - **Ex-Showroom Price: Missing**
   - **Seller Type → Dealer: 0, Individual: 1**

**Explanation**

- The original `seller_type` column (categorical: *Individual*, *Dealer*) has been replaced with two numeric columns:
  - `seller_type_Dealer` → 1 if the seller is a dealer, else 0
  - `seller_type_Individual` → 1 if the seller is an individual, else 0

**Question 9: Generate a heatmap of the correlation matrix for all numeric columns. What correlations stand out the most?**

## Answer:

```python
# Generating correlation heatmap of numeric columns from bike dataset

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

import os

# Load dataset

file_path = "/mnt/data/BIKE DETAILS.csv"

df = pd.read_csv(file_path)

# Select only numeric columns

numeric_cols = ['selling_price', 'year', 'km_driven', 'ex_showroom_price']

df_numeric = df[numeric_cols]

# Compute correlation matrix

correlation_matrix = df_numeric.corr()

# Plot heatmap

plt.figure(figsize=(8, 6))

sns.set(style="whitegrid")

sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5)

plt.title("Correlation Heatmap of Numeric Columns")

plt.tight_layout()

# Save plot

output_path = "/mnt/data/correlation_heatmap.png"

plt.savefig(output_path)

print("Correlation heatmap of numeric columns saved as correlation_heatmap.png")
```
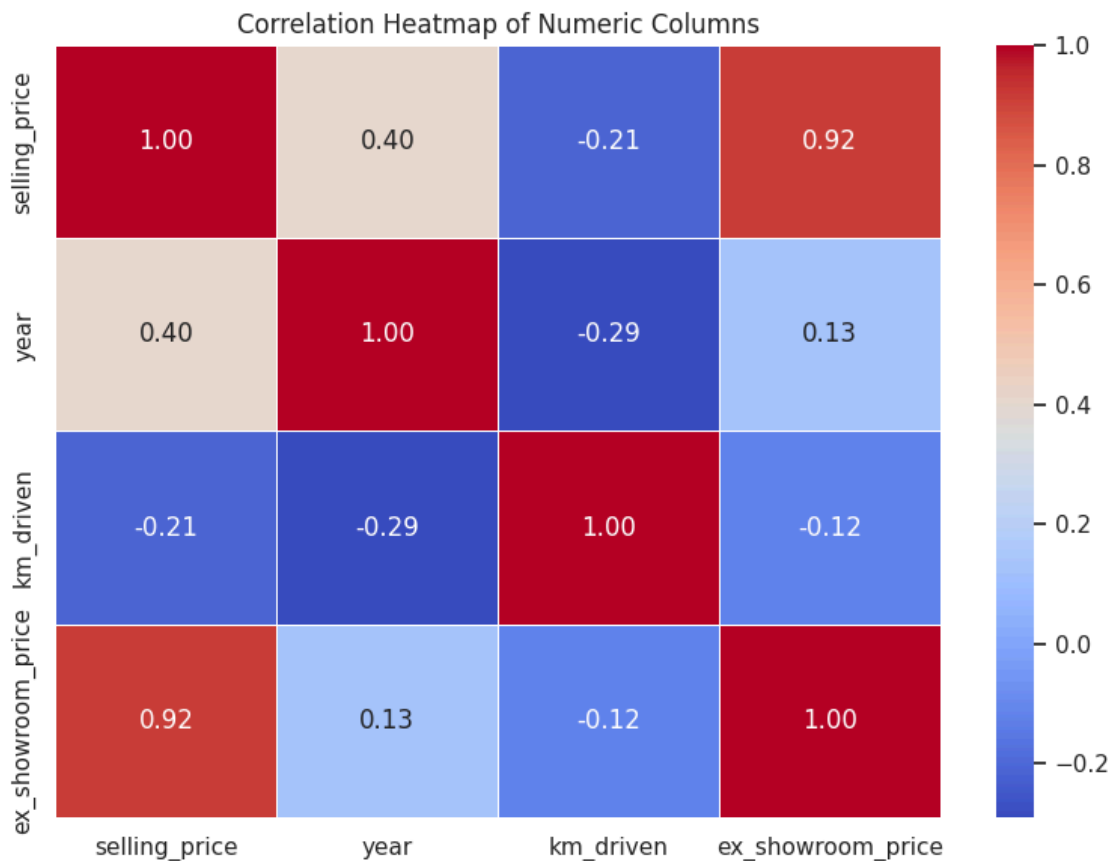
**Key Correlations That Stand Out**

- Selling Price ↔ Ex-Showroom Price → Strong positive correlation. Bikes with higher showroom prices tend to retain higher resale values.
- Year ↔ Selling Price → Positive correlation. Newer bikes generally sell for more, showing depreciation with age.
- Km Driven ↔ Selling Price → Weak negative correlation. More kilometers usually reduce resale value, but the effect is less pronounced compared to year and showroom price.

- Year ↔ Km Driven → Slight negative correlation. Older bikes have typically accumulated more kilometers.

Correlation Heatmap of Numeric Columns

|  | selling_price | year | km_driven | ex_showroom_price |
|---|---|---|---|---|
| **selling_price** | 1.00 | 0.40 | -0.21 | 0.92 |
| **year** | 0.40 | 1.00 | -0.29 | 0.13 |
| **km_driven** | -0.21 | -0.29 | 1.00 | -0.12 |
| **ex_showroom_price** | 0.92 | 0.13 | -0.12 | 1.00 |

## Question 10: Summarize your findings in a brief report:

● **What are the most important factors affecting a bike's selling price?**

● **Mention any data cleaning or feature engineering you performed.**

## Answer:

Brief Report on Bike Resale Analysis

**Key Factors Affecting Selling Price**

1. **Year of Manufacture**
   ○ Newer bikes consistently command higher resale values.
   ○ Clear depreciation trend: older bikes (2005–2012) cluster below ₹50,000, while 2018–2020 models often exceed ₹100,000.
2. **Ex-Showroom Price**

- ○ Strong positive correlation with selling price.
- ○ Premium bikes (Royal Enfield, KTM, Harley-Davidson, Kawasaki) retain higher resale values relative to their original showroom price.

3. **Kilometers Driven (Usage)**
   - ○ Negative correlation with selling price.
   - ○ More usage generally reduces resale value, though the effect is weaker compared to age and showroom price.
   - ○ Extreme mileage outliers distorted averages until cleaned.

4. **Ownership History**
   - ○ 1st owner bikes sell for more and have lower average kilometers driven.
   - ○ Resale value drops with 2nd/3rd ownership due to accumulated usage and perceived wear.

5. **Seller Type**

- ● Individual sellers tend to list bikes at higher prices than dealers, possibly reflecting better condition or premium models.

**Data Cleaning & Feature Engineering**

- ● **Missing Values**
  - ○ Dropped rows with missing `selling_price` (critical field).
  - ○ Imputed missing `ex_showroom_price` using median values grouped by bike model.
  - ○ Filled missing `km_driven` with realistic medians, capped extreme anomalies.
- ● **Outlier Treatment**
  - ○ Applied **IQR method** to `km_driven` → removed unrealistic values (e.g., 8.8 million km).
  - ○ Result: mean km dropped from ~34k to ~28k, making distribution more realistic.
- ● **Feature Engineering**
- ● One-hot encoded categorical variables (`seller_type`, `owner`) for ML readiness.
- ● Added flags for imputed values to track data quality.
- ● Explored correlations with a heatmap to identify strongest predictors of resale value.

# Overall Insight

**The most important drivers of resale price are:**

- ● **Bike's age (year)**
- ● **Original showroom price**
- ● **Ownership history**