

## **EXP.NO : 1 (setting up python environment and lib-jupyter notebook)**

**Setting up the Python environment and libraries-Jupyter Notebook: code:**

```
1.print("Hello, Jupyter Notebook!")
```

**O/P:**

→ Hello, Jupyter Notebook!

**2.import ipywidgets as widgets :**

```
from IPython.display import display  
  
slider = widgets.IntSlider(value=5, min=0, max=10, step=1, description='Number:')  
  
display(slider)  
  
button = widgets.Button(description="Click Me") display(button)  
  
def on_click(b):  
    print(f"Slider value is: {slider.value}") button.on_click(on_click)
```

**O/P:**

→ Number:  5  
Click Me  
Slider value is: 5

## **EXP.NO : 2 (EDA – data import and export):**

```
import pandas as pd  
  
from sqlalchemy import create_engine  
  
import sqlite3  
  
import requests
```

```
import matplotlib.pyplot as plt

# Load CSV

df_csv = pd.read_csv("/content/gender_submission.csv") print("CSV Data:\n",
df_csv.head())

# Export to Excel

excel_path = "titanic_full.xlsx"

df_csv.to_excel(excel_path, index=False)

print(f"Entire CSV exported to: {excel_path}")

# Read Excel file and preview

df_excel = pd.read_excel("/content/titanic_full.xlsx") print("Excel Data Preview:")

display(df_excel.head()) # ✅ Now appears in a nice box

# SQL setup and query

engine = create_engine("sqlite:///memory:")

df_csv.to_sql("titanic", engine, if_exists="replace", index=False)

query = """

SELECT Survived, COUNT(*) AS count

FROM titanic

GROUP BY Survived

"""

df_sql = pd.read_sql_query(query, engine)

print("SQL Query Result:\n", df_sql)

# Web scraping

url = "https://en.wikipedia.org/wiki/Titanic"

try:

    tables = pd.read_html(url)

    print("\nWeb Scraped Table Example:\n", tables[0].head()) except:

    print("\nUnable to scrape table from the web.")

# Plot survival count
```

```

df_sql.plot(kind='bar', x='Survived', y='count', legend=False, color=['red', 'green'])

plt.title("Titanic Survival Count")

plt.xlabel("Survived (0 = No, 1 = Yes)")

plt.ylabel("Number of Passengers")

plt.xticks(rotation=0)

plt.grid(axis='y')

plt.tight_layout()

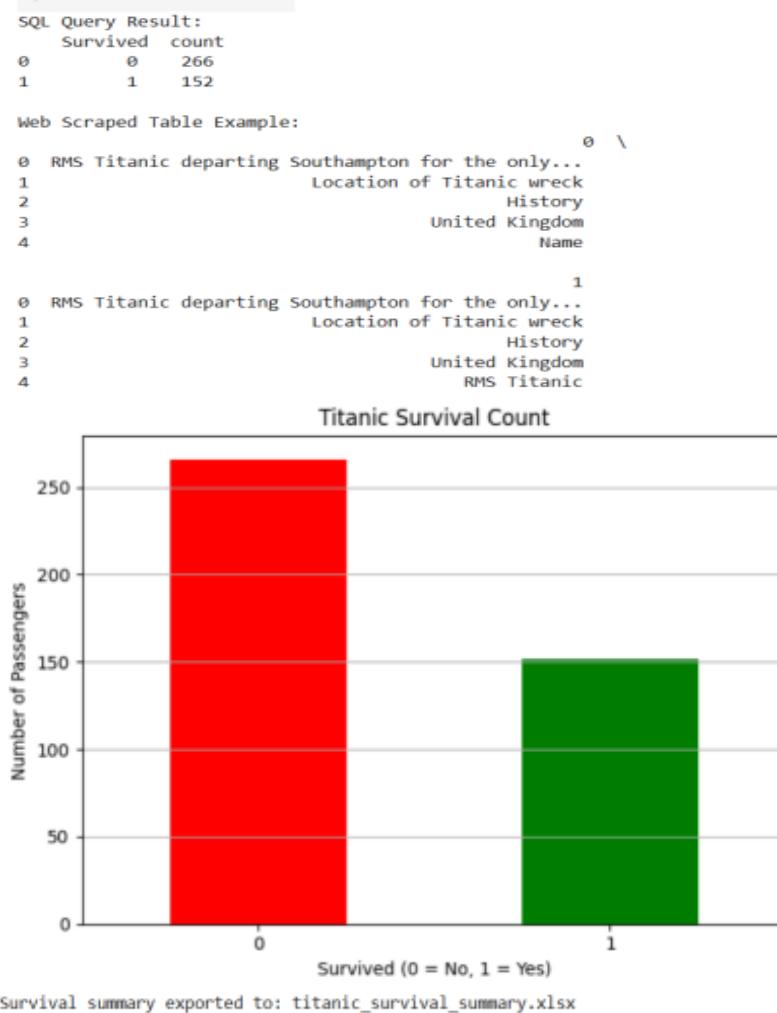
plt.show()

# Export SQL result to Excel

df_sql.to_excel("titanic_survival_summary.xlsx", index=False) print("Survival summary
exported to: titanic_survival_summary.xlsx")

```

**Output:**



## **EXP.NO : 3 (EDA – data cleaning):**

```
# Step 1: Import required libraries

import pandas as pd

import numpy as np

from sklearn.preprocessing import StandardScaler, MinMaxScaler


def main():

    # Step 2: Create a sample DataFrame

    data = {

        "Name": ["Alice", "Bob", "Charlie", np.nan, "Eve", "Alice"],

        "Age": [25, np.nan, 30, 22, 29, 25],

        "Salary": [50000, 60000, 55000, 52000, np.nan, 50000],

        "Department": ["HR", "IT", "IT", "HR", "Finance", "HR"],

    }

    df = pd.DataFrame(data)

    print("Original DataFrame:\n")

    print(df)

    # Step 3: Detect missing values

    print("\nMissing values before cleaning:")

    print(df.isnull().sum())

    # Step 4: Fill missing values using forward fill

    df = df.fillna()
```

```
# Step 5: Drop any remaining rows with missing values
df = df.dropna()

# Step 6: Show missing values after cleaning
print("\nMissing values after filling and dropping:")
print(df.isnull().sum())

# Step 7: Remove duplicate rows
print("\nDuplicate rows before removal:", df.duplicated().sum())
df = df.drop_duplicates()
print("Duplicate rows after removal:", df.duplicated().sum())

# Step 8: Remove unnecessary columns
if "Name" in df.columns:
    df = df.drop(columns=["Name"])

# Step 9: Show data types before conversion
print("\nData types before conversion:")
print(df.dtypes)

# Step 10: Convert data types for consistency
df["Age"] = df["Age"].astype(int)
df["Salary"] = df["Salary"].astype(int)

# Step 11: Show data types after conversion
print("\nData types after conversion:")
print(df.dtypes)
```

```
# Step 12: Normalize numerical columns

numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
print("\nNumeric columns to normalize:", numeric_cols)

# Standardization (Z-score)
scaler_std = StandardScaler()
df_standardized = df.copy()
df_standardized[numeric_cols] =
scaler_std.fit_transform(df_standardized[numeric_cols])

# Min-Max Scaling (0-1 range)
scaler_mm = MinMaxScaler()
df_minmax = df.copy()
df_minmax[numeric_cols] = scaler_mm.fit_transform(df_minmax[numeric_cols])

# Step 13: Save cleaned data to CSV
df.to_csv("cleaned_data.csv", index=False)
df_standardized.to_csv("standardized_data.csv", index=False)
df_minmax.to_csv("minmax_scaled_data.csv", index=False)

# Step 14: Final Outputs
print("\nFinal Cleaned DataFrame:\n")
print(df)

print("\nStandardized DataFrame:\n")
print(df_standardized)
```

```

print("\nMin-Max Scaled DataFrame:\n")

print(df_minmax)

print("\nFiles saved: cleaned_data.csv, standardized_data.csv,
minmax_scaled_data.csv")

if __name__ == "__main__":
    main()

```

**output:**

The screenshot shows a Jupyter Notebook interface with two code cells and their outputs.

**Code Cell 1 Output:**

- Data types after conversion:**

|   | Age   | Salary | Department |
|---|-------|--------|------------|
| 0 | int64 | int64  | object     |
| 1 |       |        | object     |
- Numeric columns to normalize: ['Age', 'Salary']**
- Final Cleaned DataFrame:**

|   | Age | Salary | Department |
|---|-----|--------|------------|
| 0 | 25  | 50000  | HR         |
| 1 | 25  | 60000  | IT         |
| 2 | 30  | 55000  | IT         |
| 3 | 22  | 52000  | HR         |
| 4 | 29  | 52000  | Finance    |
- Standardized DataFrame:**

|   | Age       | Salary    | Department |
|---|-----------|-----------|------------|
| 0 | -0.410152 | -1.089725 | HR         |
| 1 | -0.410152 | 1.777972  | IT         |
| 2 | 1.298813  | 0.344124  | IT         |
| 3 | -1.435530 | -0.516185 | HR         |
| 4 | 0.057020  | -0.516185 | Finance    |
- Min-Max scaled DataFrame:**

|   | Age   | Salary | Department |
|---|-------|--------|------------|
| 0 | 0.375 | 0.0    | HR         |
| 1 | 0.375 | 1.0    | IT         |
| 2 | 1.000 | 0.5    | IT         |
| 3 | 0.800 | 0.2    | HR         |
| 4 | 0.875 | 0.2    | Finance    |

**Code Cell 2 Output:**

- Original DataFrame:**

|   | Name    | Age  | Salary  | Department |
|---|---------|------|---------|------------|
| 0 | Alice   | 25.0 | 50000.0 | HR         |
| 1 | Bob     | NaN  | 60000.0 | IT         |
| 2 | Charlie | 30.0 | 55000.0 | IT         |
| 3 | NaN     | 22.0 | 52000.0 | HR         |
| 4 | Eve     | 29.0 | NaN     | Finance    |
| 5 | Alice   | 25.0 | 50000.0 | HR         |
- Missing values before cleaning:**

|   | Name | Age | Salary | Department |
|---|------|-----|--------|------------|
| 0 | 1    | 1   | 1      | 0          |
- Missing values after filling and dropping:**

|   | Name | Age | Salary | Department |
|---|------|-----|--------|------------|
| 0 | 0    | 0   | 0      | 0          |
- Duplicate rows before removal: 1**
- Duplicate rows after removal: 0**
- Data types before conversion:**

|   | Age     | Salary  | Department |
|---|---------|---------|------------|
| 0 | float64 | float64 | object     |
| 1 |         |         | object     |

## **EXP.NO : 4 (EDA – data inspection and analysis):**

```
# Step 1: Import necessary libraries
```

```
import pandas as pd
```

```
import numpy as np
```

```
# Step 2: Create a sample DataFrame
```

```
data = {
```

```
    'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
```

```
    'Age': [25, 30, 22, 28, 35],
```

```
    'Salary': [50000, 60000, 52000, 58000, 65000],
```

```
    'Department': ['HR', 'IT', 'HR', 'Finance', 'IT']
```

```
}
```

```
df = pd.DataFrame(data)
```

```
# Step 3: View and inspect the DataFrame
```

```
print("First 5 rows of the DataFrame:\n")
```

```
print(df.head())
```

```
print("\nInfo about DataFrame:")
```

```
print(df.info())
```

```
print("\nShape of DataFrame:", df.shape)
```

```
print("Column Names:", df.columns.tolist())
```

```
# Step 4: Filter and subset data using conditions
```

```
print("\nEmployees with Salary > 55000:")
```

```
high_salary = df[df['Salary'] > 55000]
```

```
print(high_salary)
```

```
print("\nEmployees from IT Department:")
it_employees = df[df['Department'] == 'IT']
print(it_employees)

# Step 5: Descriptive statistics - Central Tendency
print("\nMean Age:", df['Age'].mean())
print("Median Age:", df['Age'].median())
print("Mode Age:", df['Age'].mode()[0])

# Step 6: Descriptive statistics - Dispersion
range_age = df['Age'].max() - df['Age'].min()
print("\nRange of Age:", range_age)
print("Variance of Age:", df['Age'].var())
print("Standard Deviation of Age:", df['Age'].std())

range_salary = df['Salary'].max() - df['Salary'].min()
print("\nRange of Salary:", range_salary)
print("Variance of Salary:", df['Salary'].var())
print("Standard Deviation of Salary:", df['Salary'].std())
```

## output:

```
22  ✓ Shape of DataFrame: (5, 4)
    ✓ Column Names: ['Name', 'Age', 'Salary', 'Department']
    ★ Employees with Salary > 55000:
      Name  Age  Salary  Department
    1  Bob   30   60000        IT
    3 David  28   58000    Finance
    4 Eve    35   65000        IT

    ★ Employees from IT Department:
      Name  Age  Salary  Department
    1  Bob   30   60000        IT
    4 Eve    35   65000        IT

    ⚡ Mean Age: 28.0
    ⚡ Median Age: 28.0
    ⚡ Mode Age: 22

    ✅ Range of Age: 13
    ✅ Variance of Age: 24.5
    ✅ Standard Deviation of Age: 4.949747468305833

    ✅ Range of Salary: 15000
    ✅ Variance of Salary: 37000000.0
    ✅ Standard Deviation of Salary: 6082.76253029822
    🔎 First 5 rows of the DataFrame:
    ↴
      Name  Age  Salary  Department
    0  Alice  25   50000        HR
    1    Bob  38   68000        IT
    2 Charlie  22   52000        HR
    3  David  28   58000    Finance
    4    Eve  35   65000        IT

    📈 Info about DataFrame:
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 5 entries, 0 to 4
    Data columns (total 4 columns):
     #   Column       Non-Null Count  Dtype  
    --- 
     0   Name         5 non-null      object 
     1   Age          5 non-null      int64  
     2   Salary        5 non-null      int64  
     3   Department    5 non-null      object 
    dtypes: int64(2), object(2)
    memory usage: 292.0+ bytes
    None

    ✓ Shape of DataFrame: (5, 4)
    ✓ Column Names: ['Name', 'Age', 'Salary', 'Department']

    ★ Employees with Salary > 55000:
      Name  Age  Salary  Department
    1  Bob   30   60000        IT
    3 David  28   58000    Finance
    4 Eve    35   65000        IT
```

## **EXP.NO : 5 (EDA – data visualization with matplotlib):**

```
# Step 1: Import libraries
import pandas as pd
import matplotlib.pyplot as plt

# Step 2: Create a sample DataFrame
data = {
    'PassengerId': [1, 2, 3, 4, 5, 6],
    'Survived': [0, 1, 1, 0, 1, 0],
    'Name': ['Allen', 'Braund', 'Cumings', 'Dawson', 'Evans', 'Frank'],
    'Age': [22, 38, 26, 29, 35, 42],
}
df = pd.DataFrame(data)

# Step 3: Line chart – Age vs PassengerId
plt.figure(figsize=(6, 4))
plt.plot(df['PassengerId'], df['Age'], marker='o')
plt.title('Line Chart: Age vs PassengerId')
plt.xlabel('PassengerId')
plt.ylabel('Age')
plt.grid(True)
plt.show()

# Step 4: Bar chart – Survived counts
survival_counts = df['Survived'].value_counts()
plt.figure(figsize=(5, 4))
plt.bar(survival_counts.index, survival_counts.values, color=['red', 'green'])
plt.title('Bar Chart: Survival Count')
plt.xlabel('Survived (0 = No, 1 = Yes)')
plt.ylabel('Count')
```

```

plt.xticks([0, 1])

plt.show()

# Step 5: Histogram – Age distribution

plt.figure(figsize=(6, 4))

plt.hist(df['Age'], bins=5, color='skyblue', edgecolor='black')

plt.title('Histogram: Age Distribution')

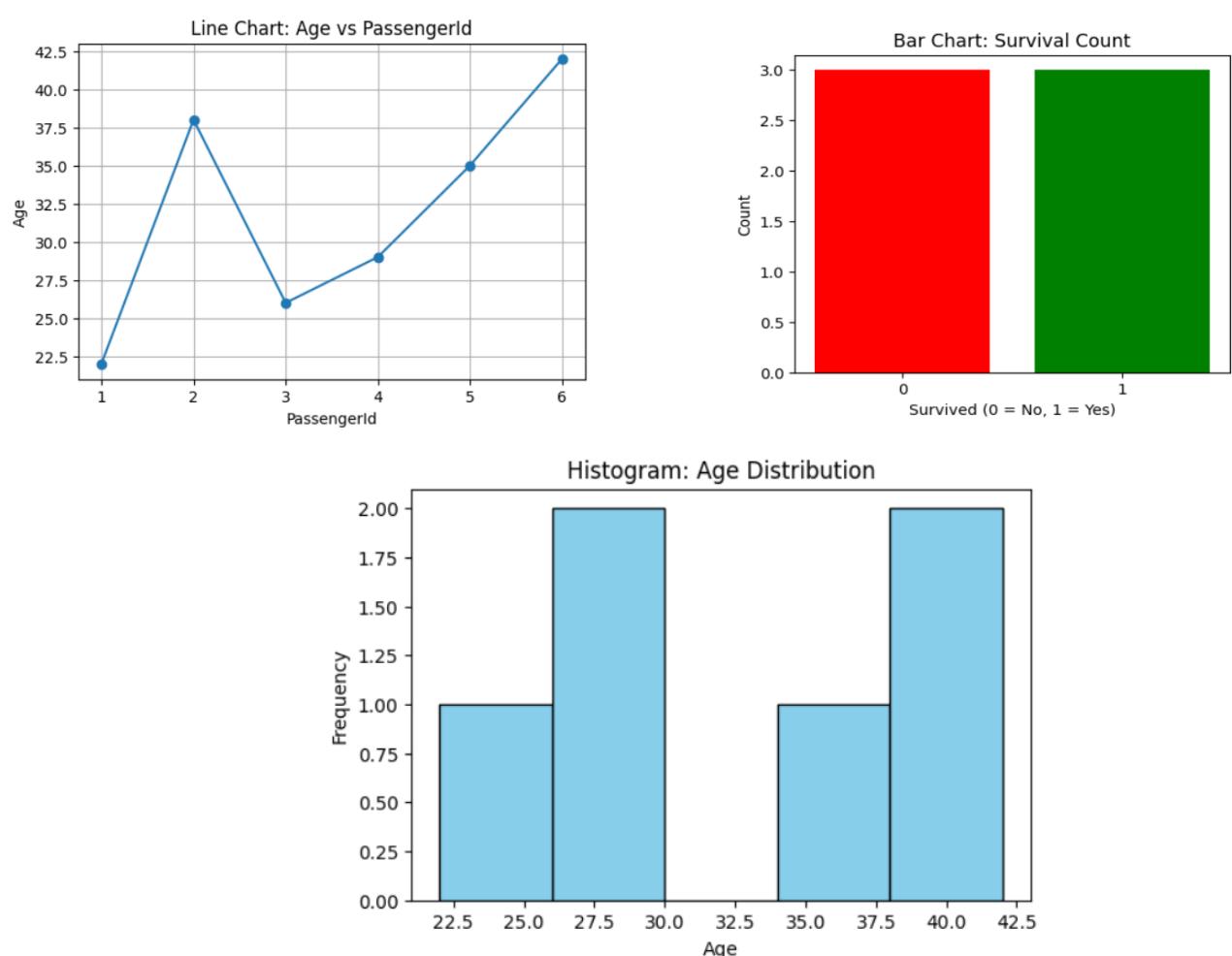
plt.xlabel('Age')

plt.ylabel('Frequency')

plt.show()

```

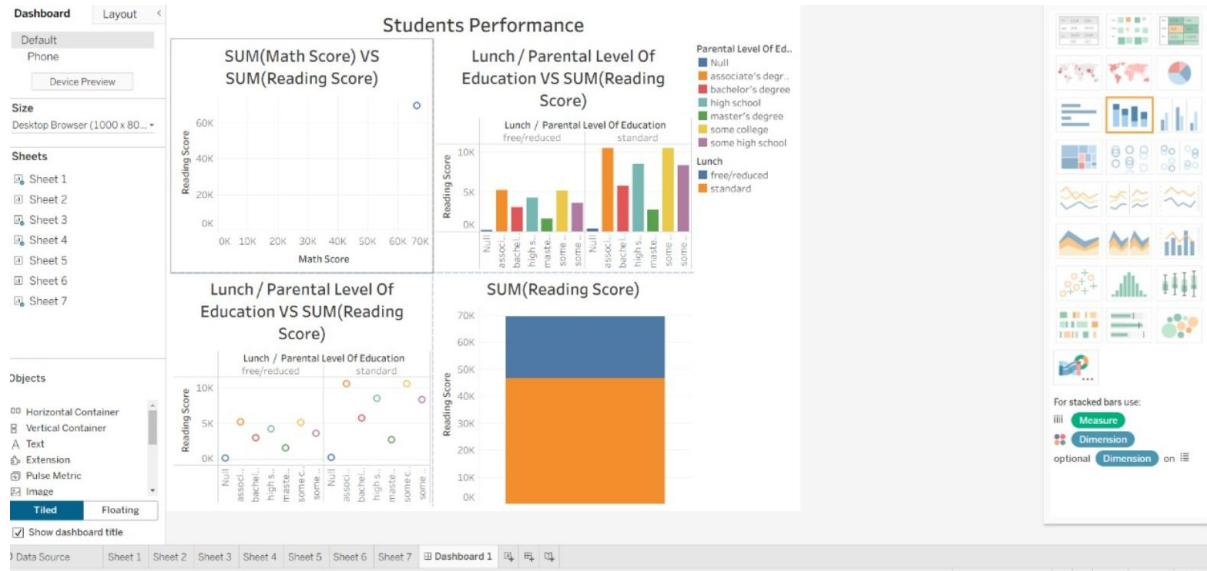
**output:**



## EXP.NO : 6:



## EXP.NO : 7:



## EXP NO 8 :

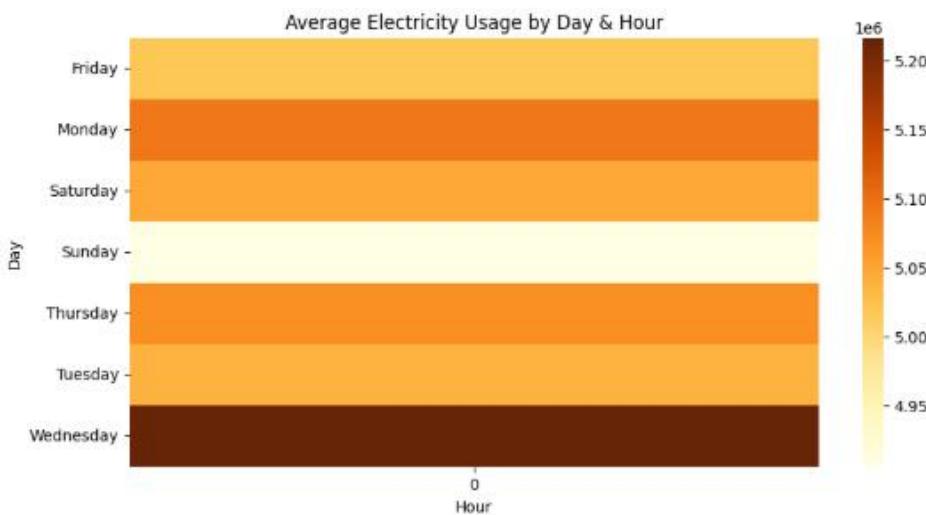
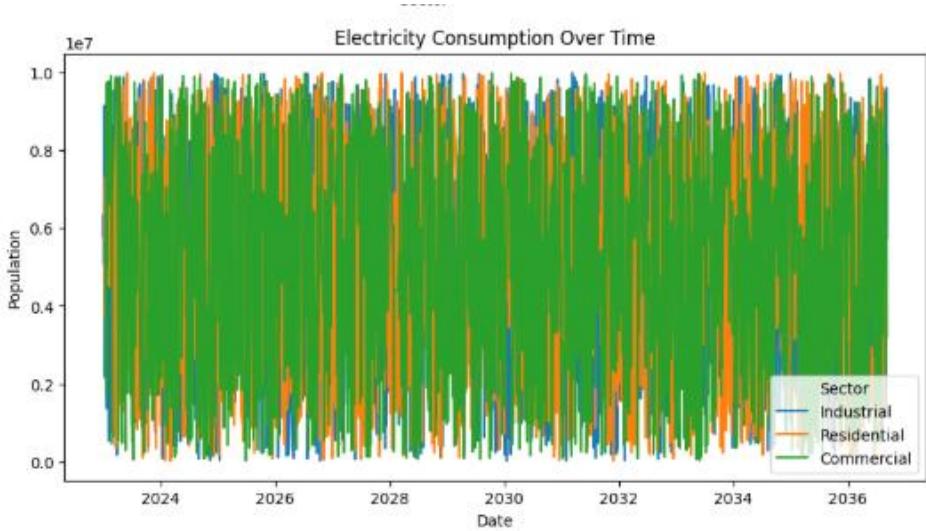
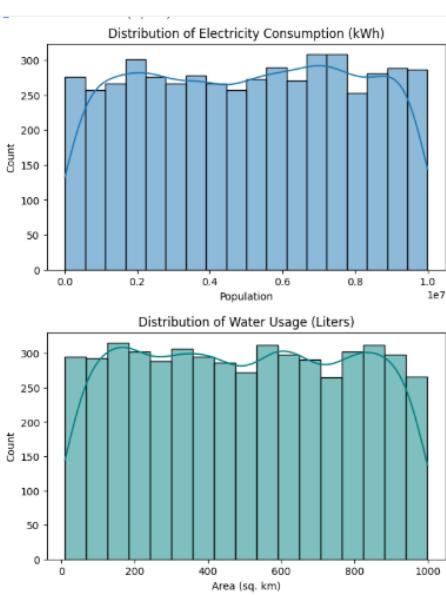
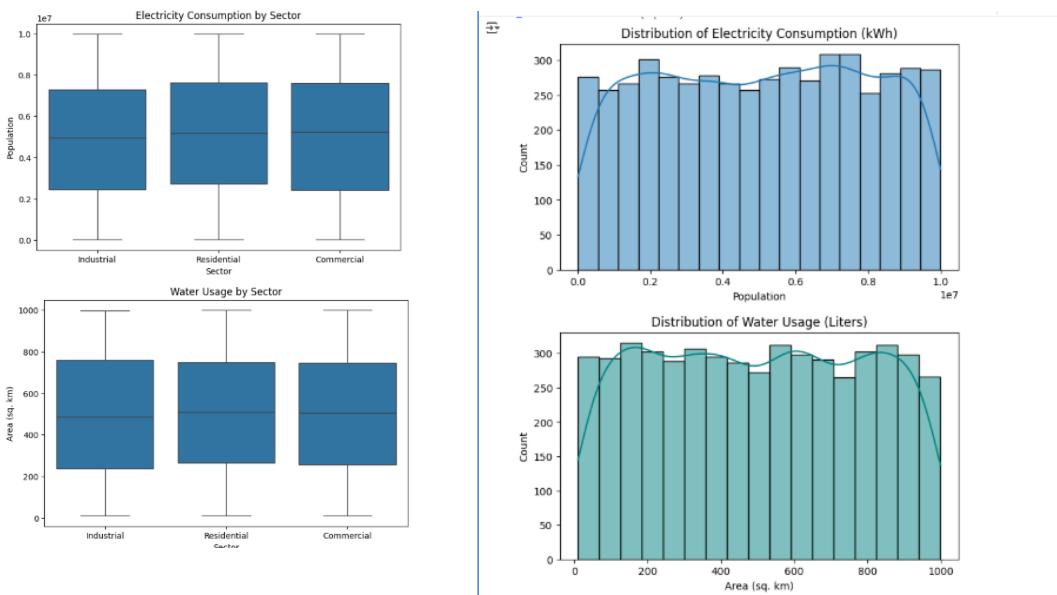
```
[5] ✓ 2s
 1
 2 # ✨ Step 1: Import libraries
 3 import pandas as pd
 4 import numpy as np
 5 import matplotlib.pyplot as plt
 6 import seaborn as sns
 7
 8 # -----
 9 # ✨ Step 2: Load the dataset
10 df = pd.read_csv("/content/smart_city_dataset (1).csv")
11
12 # Check columns and basic info
13 print("✅ Columns in the dataset:")
14 print(df.columns.tolist())
15 print("\nDataset shape:", df.shape)
16 print("\nMissing values:\n", df.isna().sum())
17
18 # -----
19 # ✨ Step 3: Identify the correct date/time column
20 # Replace below with the actual date/time column name in your dataset
21 # (example: 'timestamp', 'date_time', 'recorded_at', etc.)
22 date_col = None
23 for col in df.columns:
24     if 'date' in col.lower() or 'time' in col.lower():
25         date_col = col
26         break
27
28 if date_col:
29     print(f"\n⌚ Using '{date_col}' as Date column.")
30     df['Date'] = pd.to_datetime(df[date_col], errors='coerce')
31 else:
32     print("\n⚠️ No date/time column found! Creating synthetic date column.")
33     df['Date'] = pd.date_range(start='2023-01-01', periods=len(df), freq='D')
34
35 # -----
36 # ✨ Step 4: Basic data cleaning
37 # Remove duplicates
38 df = df.drop_duplicates()
39
40 # Fill numeric columns with median, categorical with mode
41 for col in df.select_dtypes(include=[np.number]).columns:
42     df[col] = df[col].fillna(df[col].median())
43
44 for col in df.select_dtypes(exclude=[np.number]).columns:
45     df[col] = df[col].fillna(df[col].mode()[0])
46
47 print("\n✅ Missing values after cleaning:", df.isna().sum().sum())
48
49 # -----
50 # ✨ Step 5: Feature engineering
51 df['Hour'] = df['Date'].dt.hour if df['Date'].dt.hour.notna().any() else np.nan
52 df['Day'] = df['Date'].dt.day_name()
53
54 # If sector info missing, assume "Residential" for demonstration
55 if 'Sector' not in df.columns and 'sector' not in df.columns:
56     df['Sector'] = np.random.choice(['Residential', 'Commercial', 'Industrial'], len(df))
57
```

```

58 # Rename likely consumption columns (auto-detect)
59 for col in df.columns:
60     if 'electric' in col.lower():
61         elec_col = col
62     if 'water' in col.lower():
63         water_col = col
64
65 # Fallbacks if not detected
66 elec_col = locals().get('elec_col', df.select_dtypes(include=[np.number]).columns[0])
67 water_col = locals().get('water_col', df.select_dtypes(include=[np.number]).columns[1])
68
69 print(f"\n⚡ Electricity Column: {elec_col}")
70 print(f"💧 Water Column: {water_col}")
71
72 # -----
73 # ✨ Step 6: Univariate Analysis
74
75 plt.figure(figsize=(7,4))
76 sns.histplot(df[elec_col], kde=True)
77 plt.title("Distribution of Electricity Consumption (kWh)")
78 plt.show()
79
80 plt.figure(figsize=(7,4))
81 sns.histplot(df[water_col], kde=True, color='teal')
82 plt.title("Distribution of Water Usage (Liters)")
83 plt.show()
84

85 # -----
86 # ✨ Step 7: Bivariate Analysis
87
88 plt.figure(figsize=(8,5))
89 sns.boxplot(x='Sector', y=elec_col, data=df)
90 plt.title("Electricity Consumption by Sector")
91 plt.show()
92
93 plt.figure(figsize=(8,5))
94 sns.boxplot(x='Sector', y=water_col, data=df)
95 plt.title("Water Usage by Sector")
96 plt.show()
97
98 # -----
99 # ✨ Step 8: Time-based Trend
100 plt.figure(figsize=(10,5))
101 sns.lineplot(x='Date', y=elec_col, hue='Sector', data=df)
102 plt.title("Electricity Consumption Over Time")
103 plt.show()
104
105 # -----
106 # ✨ Step 9: Heatmap - Peak Usage Hours
107 if df['Hour'].notna().any():
108     pivot_data = df.pivot_table(values=elec_col, index='Day', columns='Hour', aggfunc='mean')
109     plt.figure(figsize=(10,5))
110     sns.heatmap(pivot_data, cmap='YlOrBr')
111     plt.title("Average Electricity Usage by Day & Hour")
112     plt.show()
113
114 # -----
115 # ✨ Step 10: Save cleaned dataset
116 df.to_csv("cleaned_smart_city_dataset.csv", index=False)
117 print("\n💾 Cleaned dataset saved as 'cleaned_smart_city_dataset.csv'")
118
119 # -----

```



## Quick summary

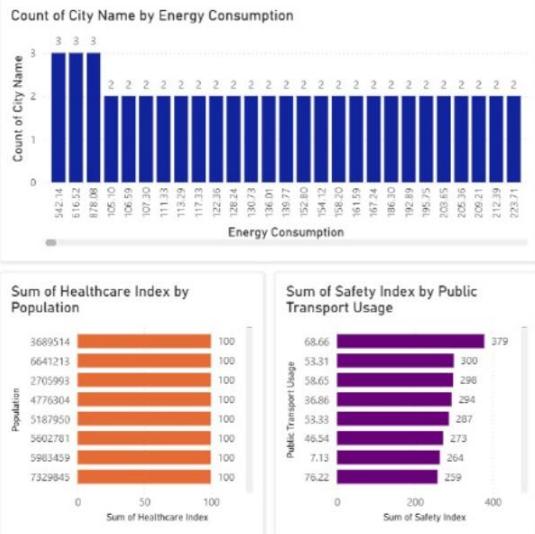
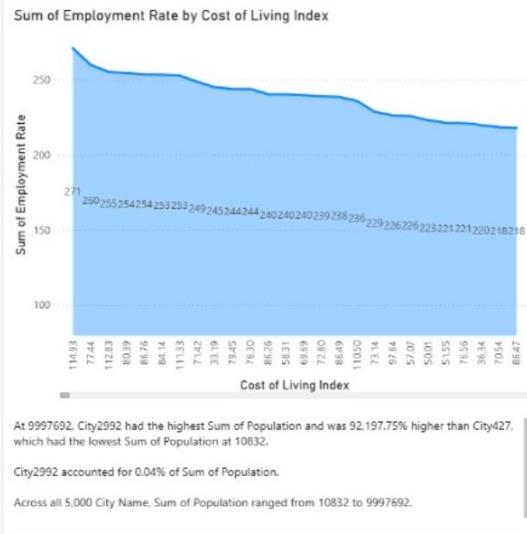
Table

25276479139  
Sum of Population

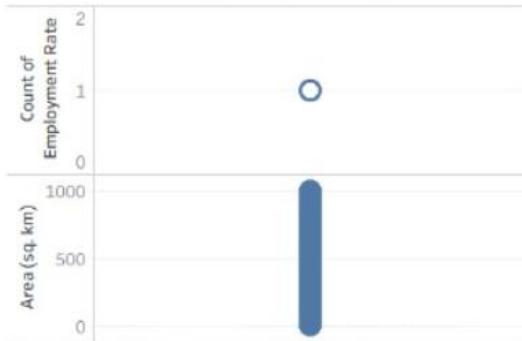
2,480,903.51  
Sum of Internet Speed (...)

2,505,485.91  
Sum of Area (sq. km)

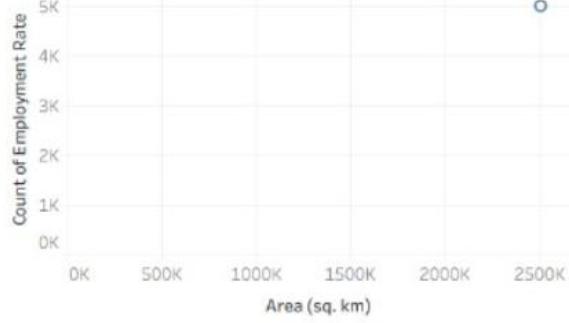
2,736,715.70  
Sum of Energy Consump...



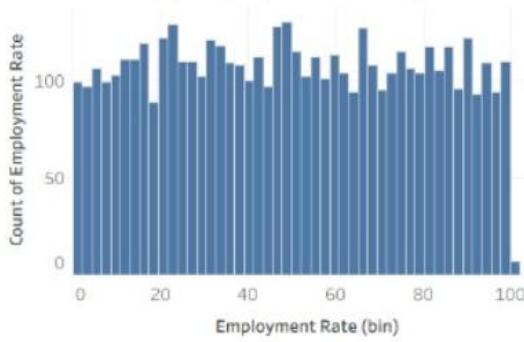
CNT(Employement Rate) /  
SUM(Area(sq.km))



SUM(Area(sq.km)) VS CNT(Employement Rate)



Employment Rate(bin) VS  
CNT(Employement Rate)



CNT(Employement Rate) /  
SUM(Area(sq.km)) VS City Name

