

localhost:8889/notebooks/DatavisualizationChap6n7.ipynb

jupyter DatavisualizationChap6n7 Last Checkpoint: 19 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38)

CHAPTER 6

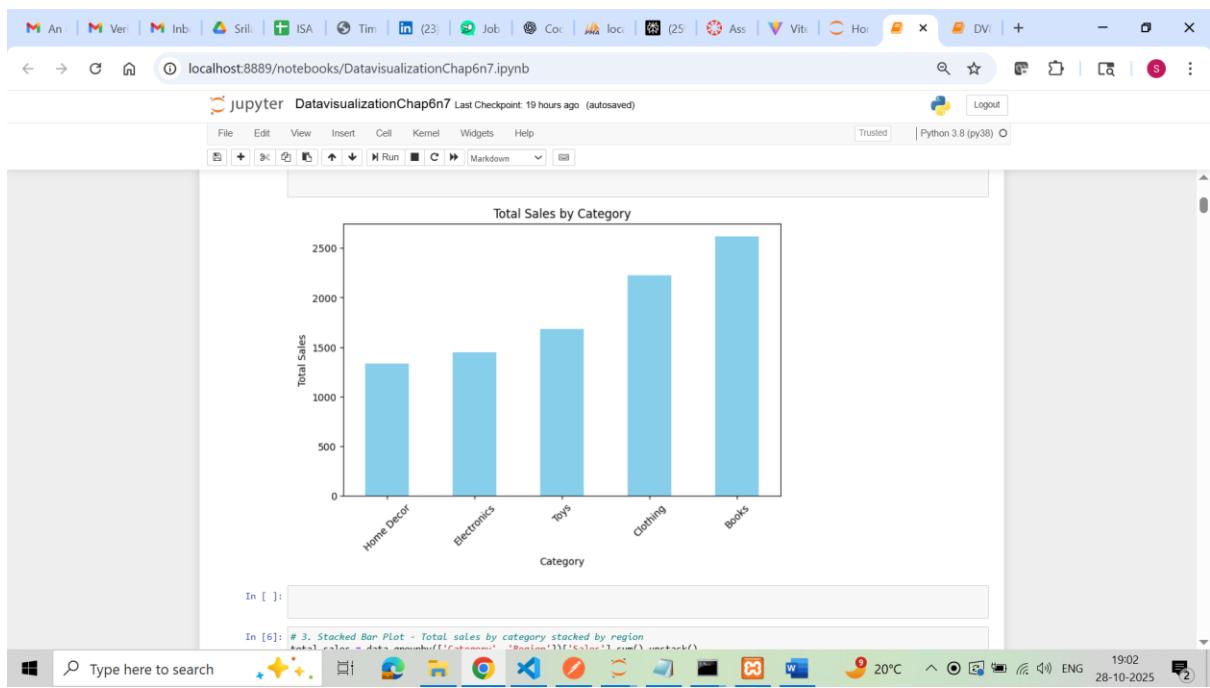
```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

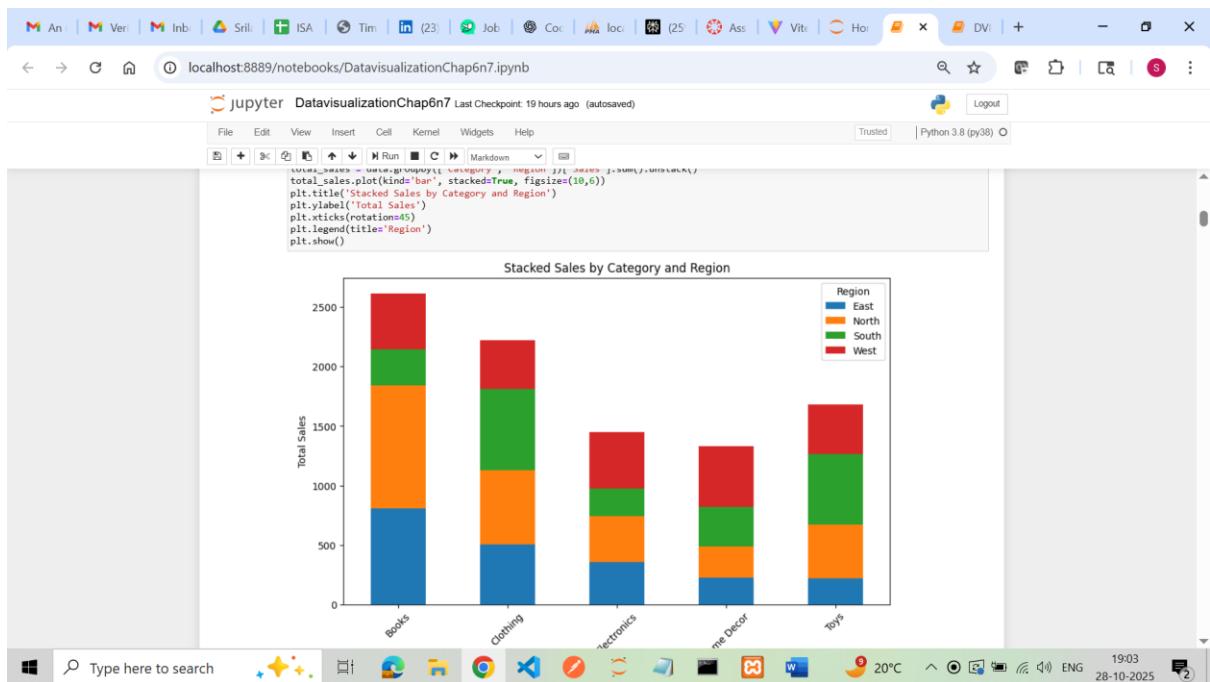
# Create synthetic dataset with 100 records
np.random.seed(42)
data = pd.DataFrame({
    'Category': np.random.choice(['Electronics', 'Clothing', 'Home Decor', 'Books', 'Toys'], 100),
    'Region': np.random.choice(['West', 'East', 'North', 'South'], 100),
    'Sales': np.random.randint(20, 150, 100)
})
```

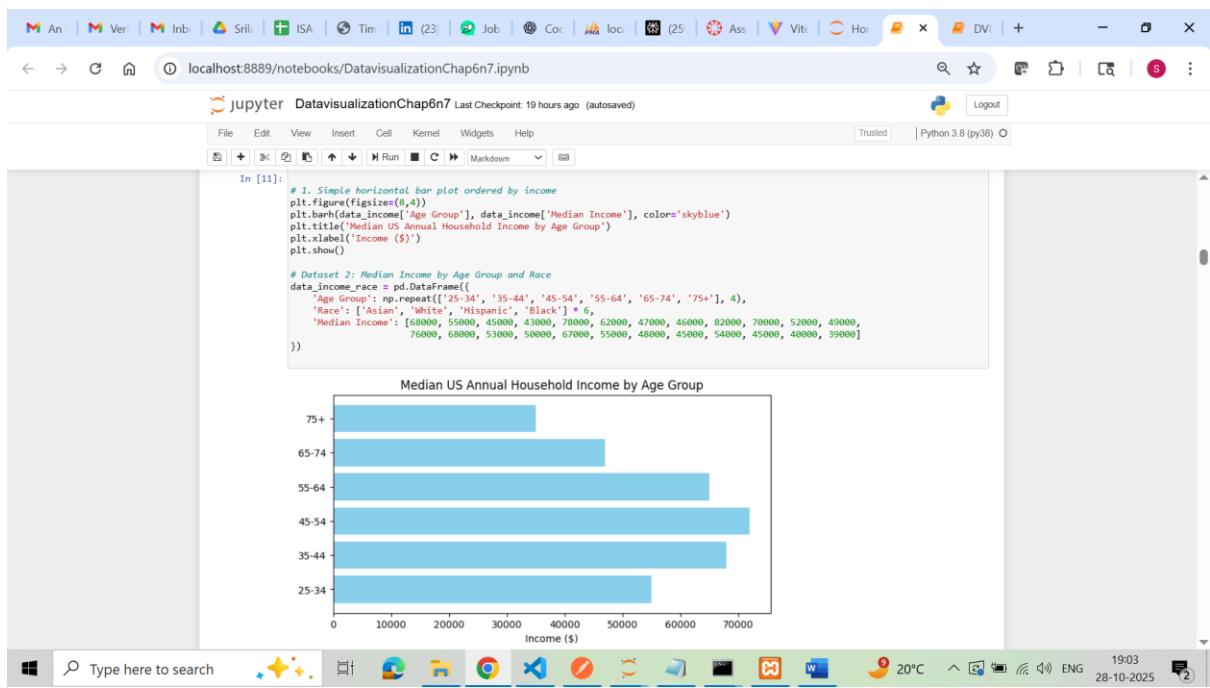
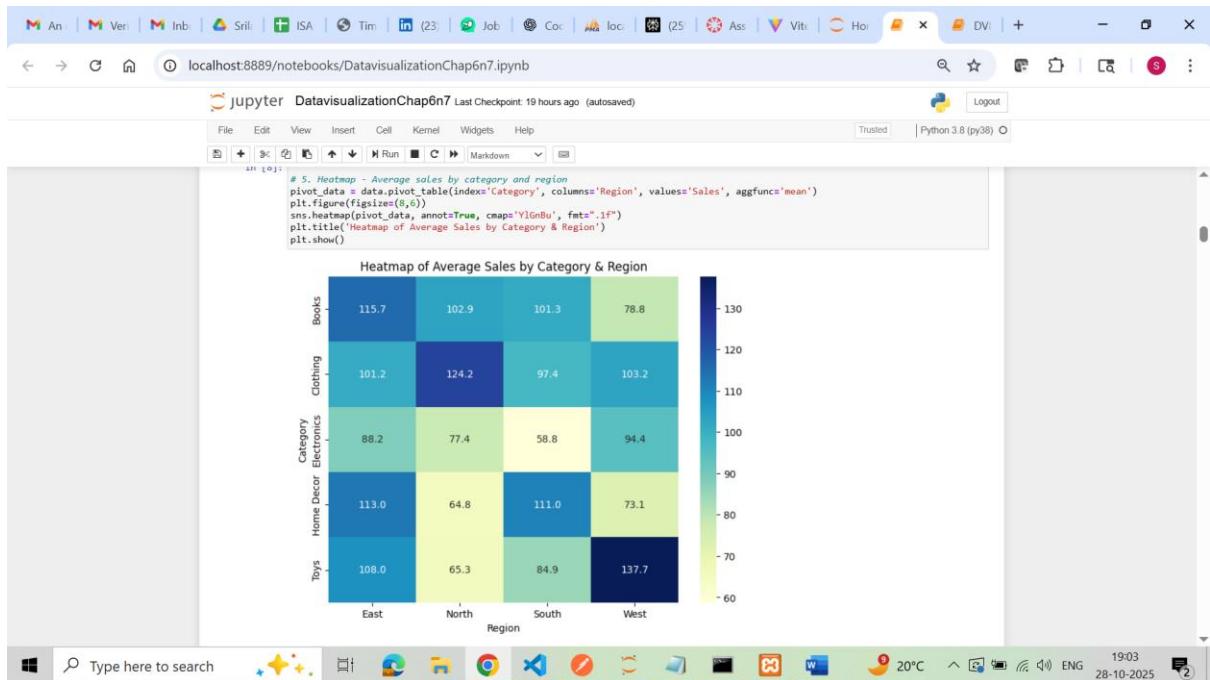
```
In [4]: # 1. Bar Plot - Total sales by category
plt.figure(figsize=(8,5))
total_sales_by_cat = data.groupby('Category')['Sales'].sum().sort_values()
total_sales_by_cat.plot(kind='bar', color='skyblue')
plt.title('Total Sales by Category')
plt.xlabel('Total Sales')
plt.xticks(rotation=45)
plt.show()
```

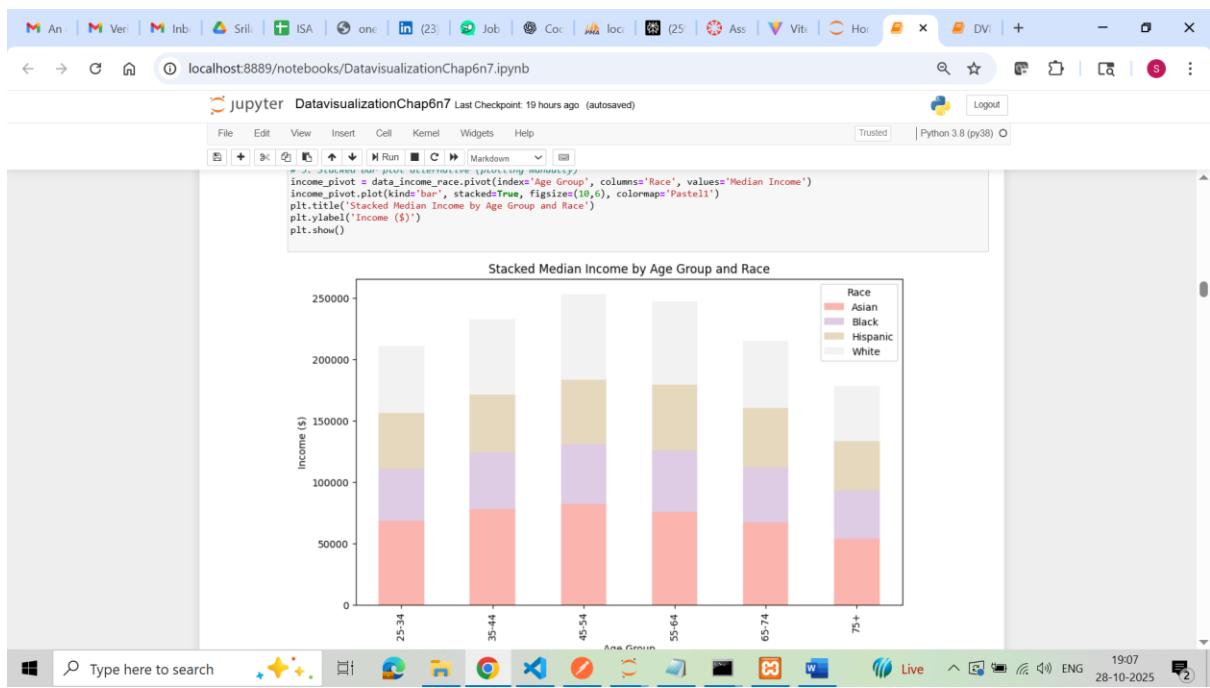
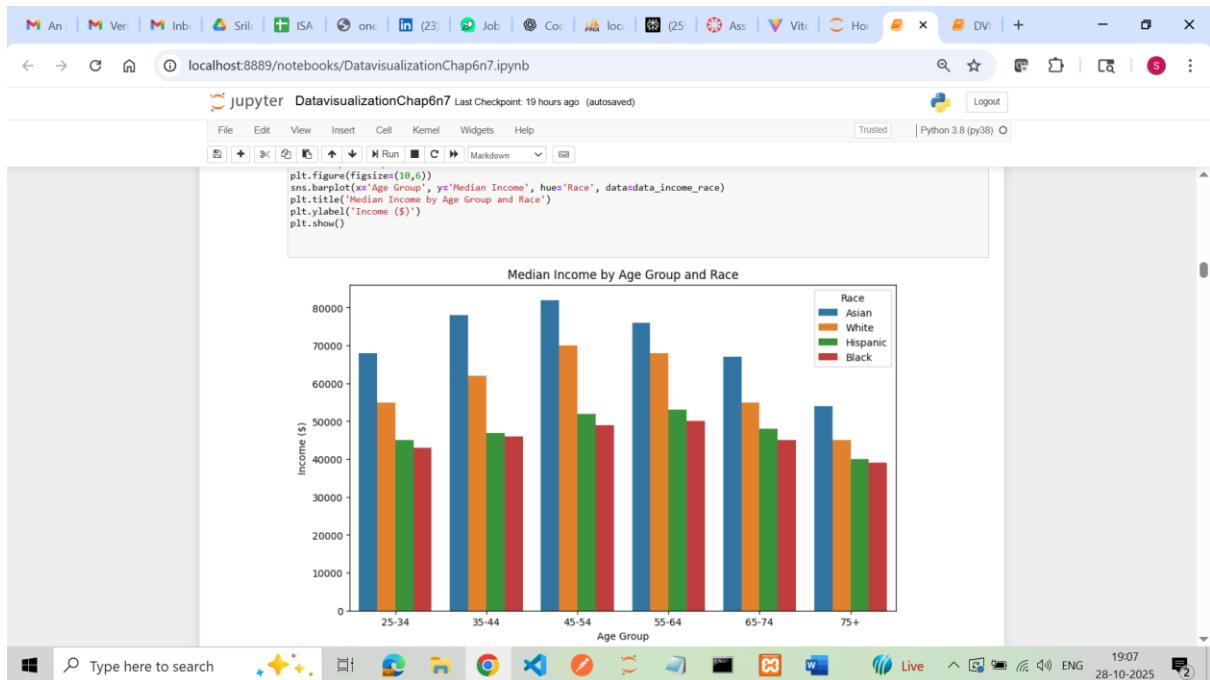
Total Sales by Category

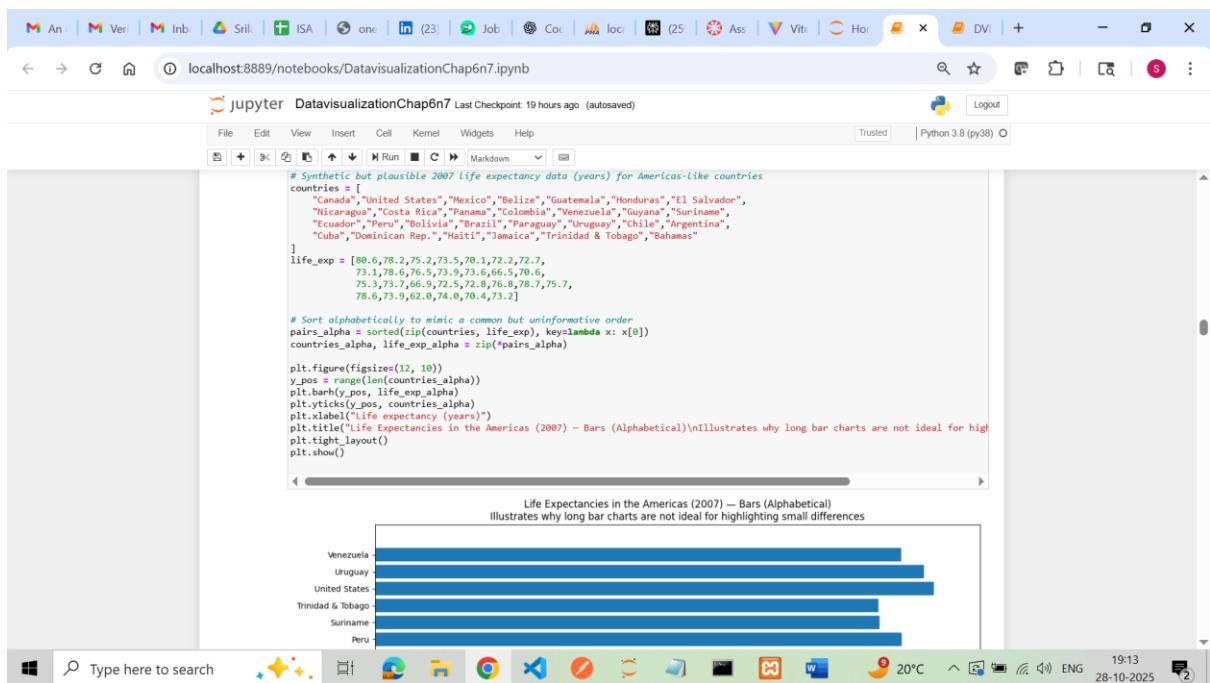
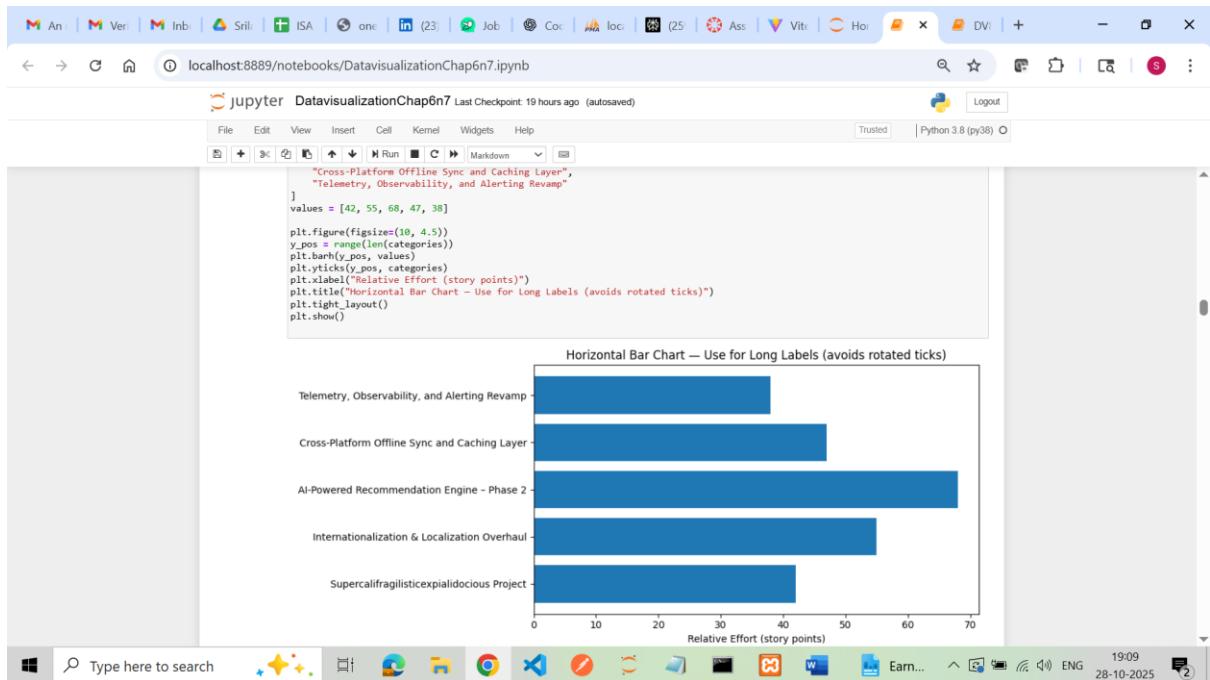
Category	Total Sales
Home Decor	1400
Electronics	1500
Toys	1700
Clothing	2200
Books	2600

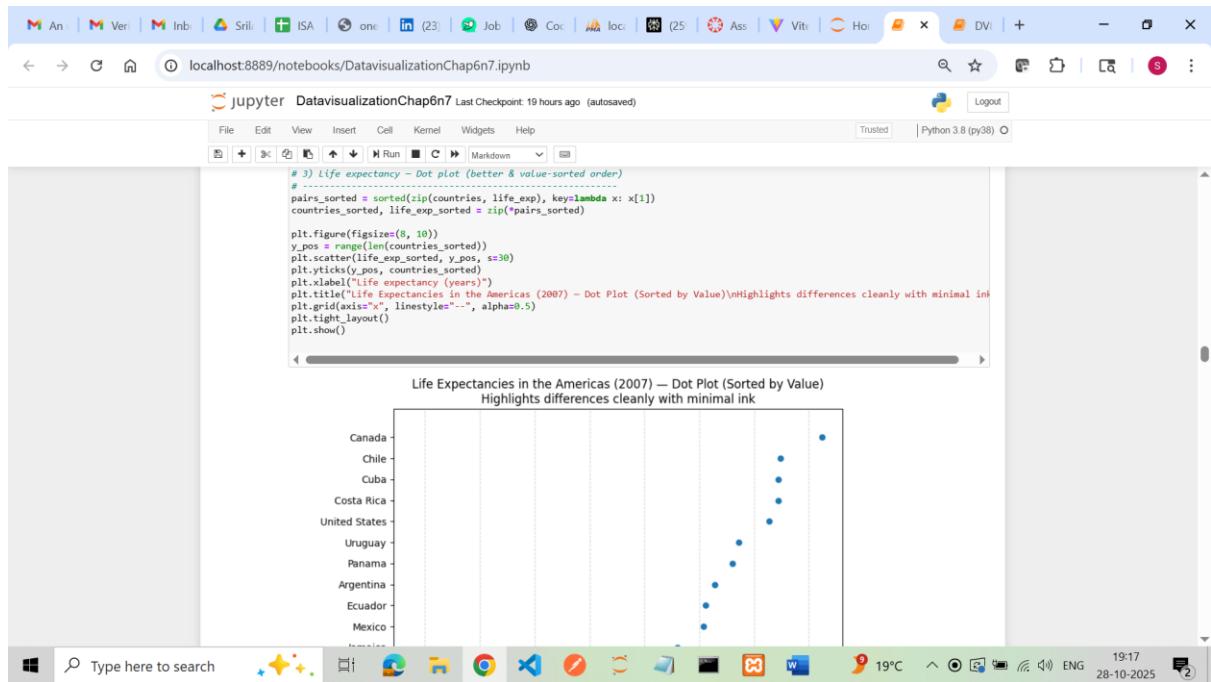
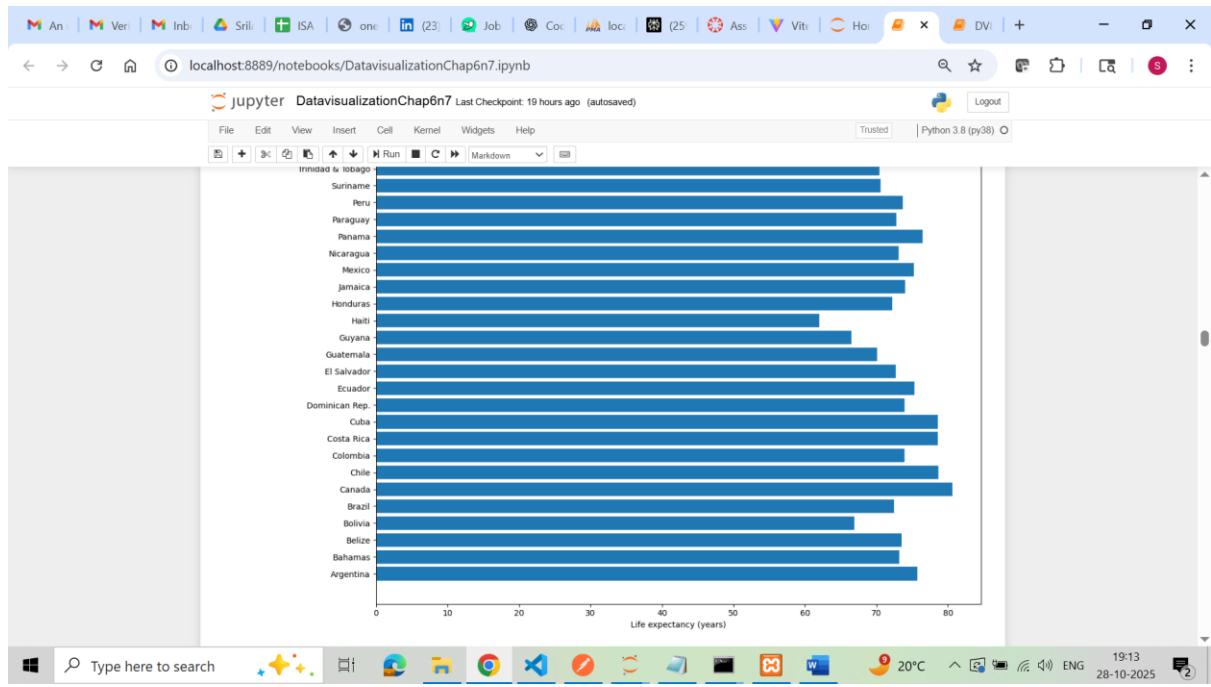


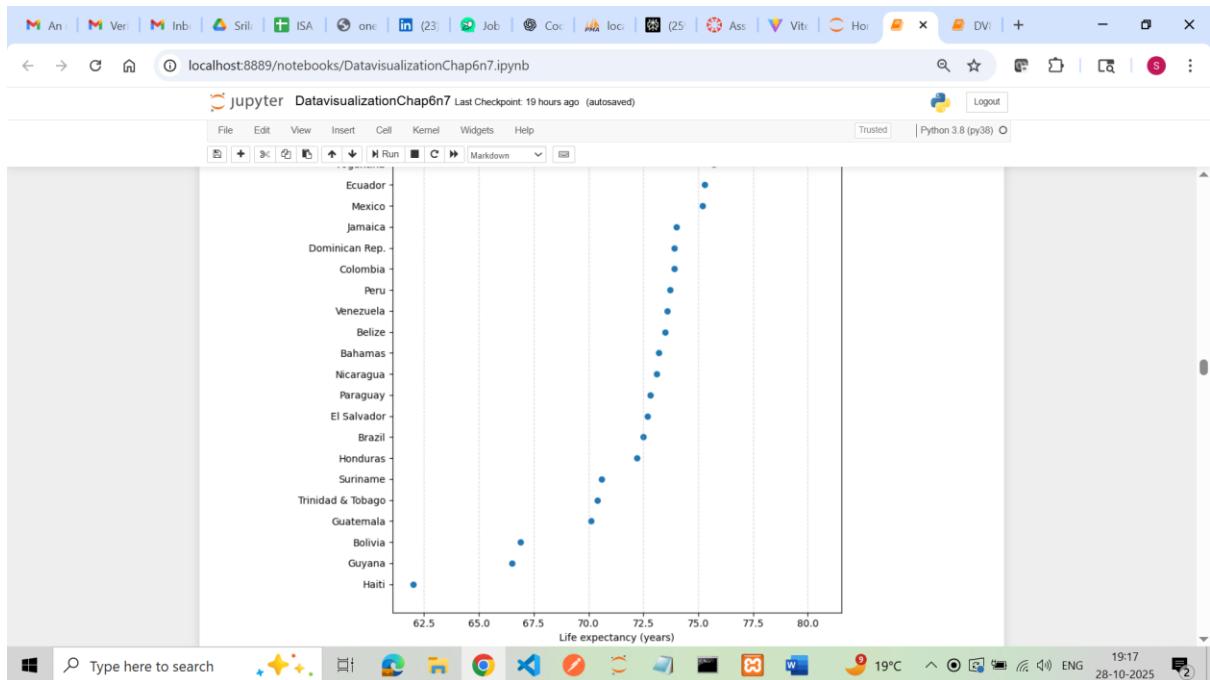




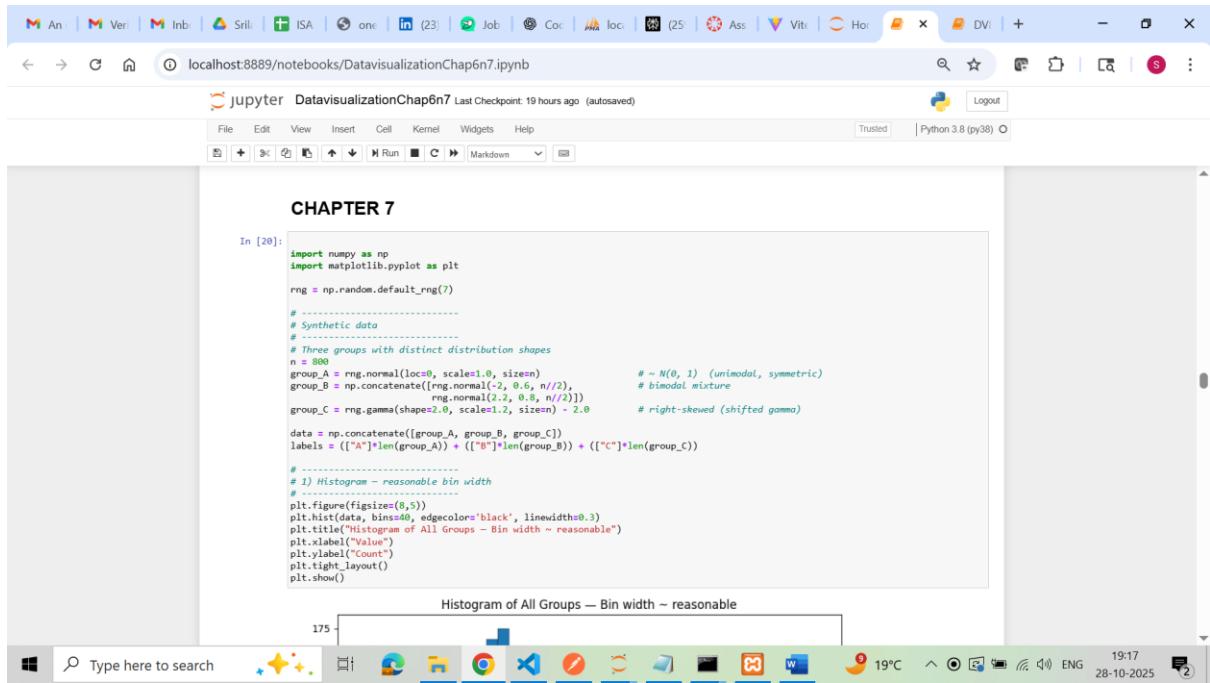


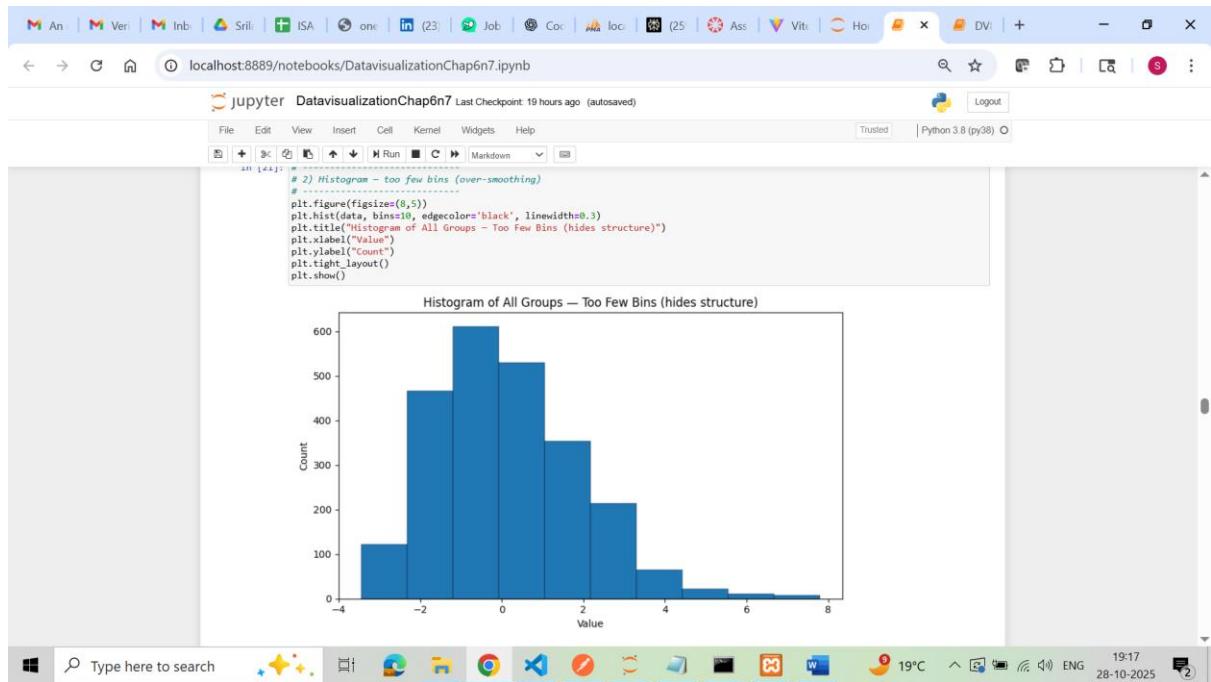
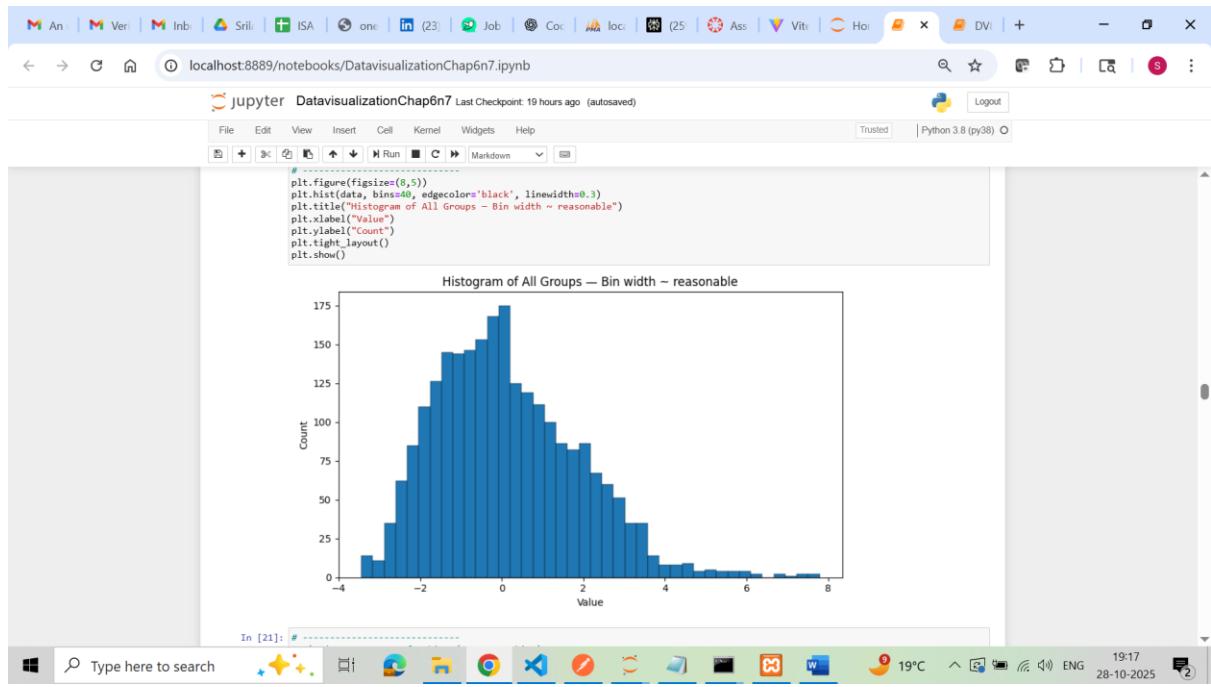


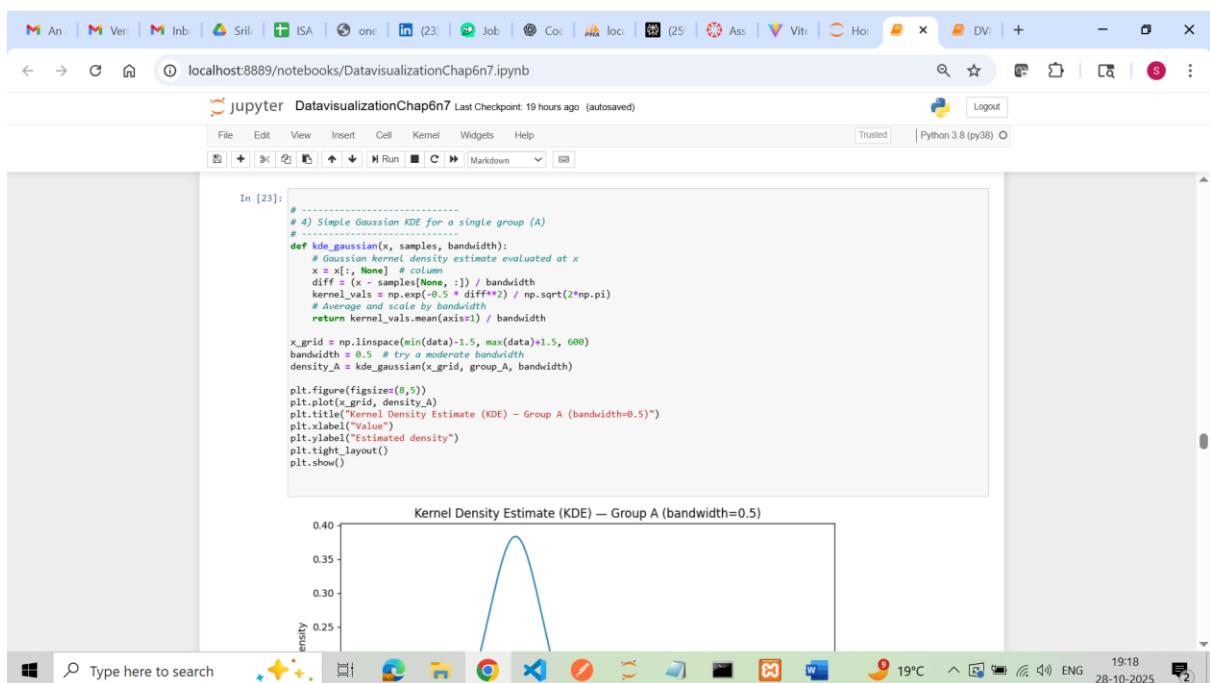
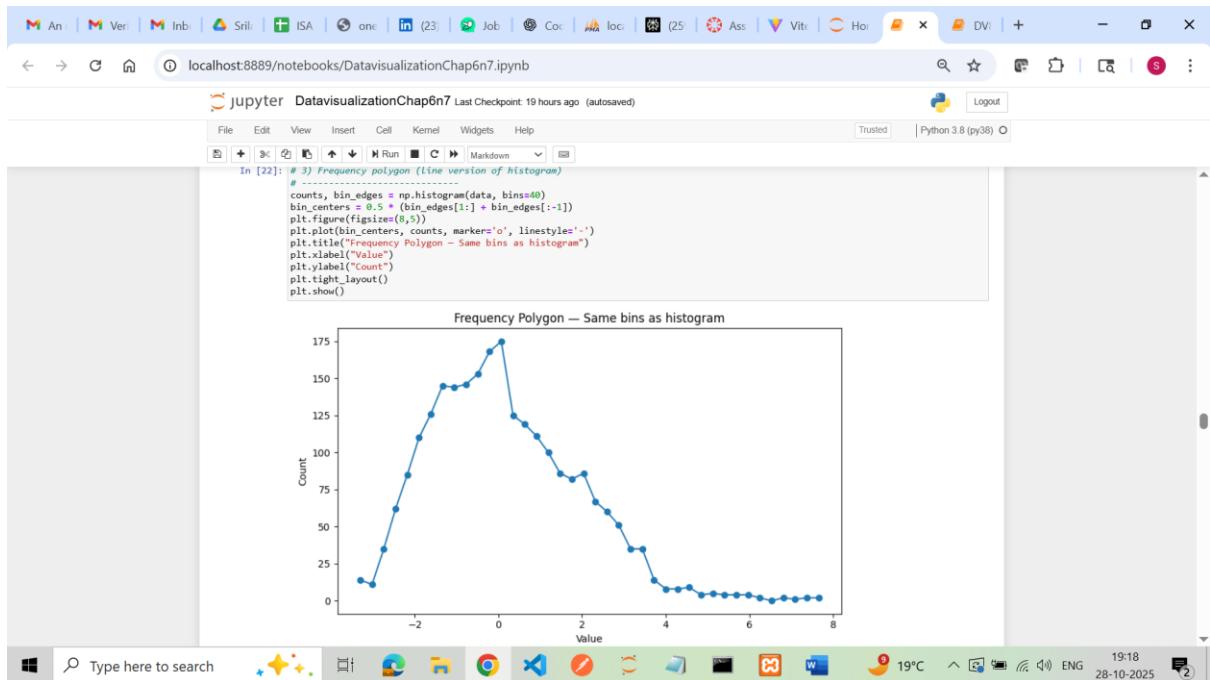


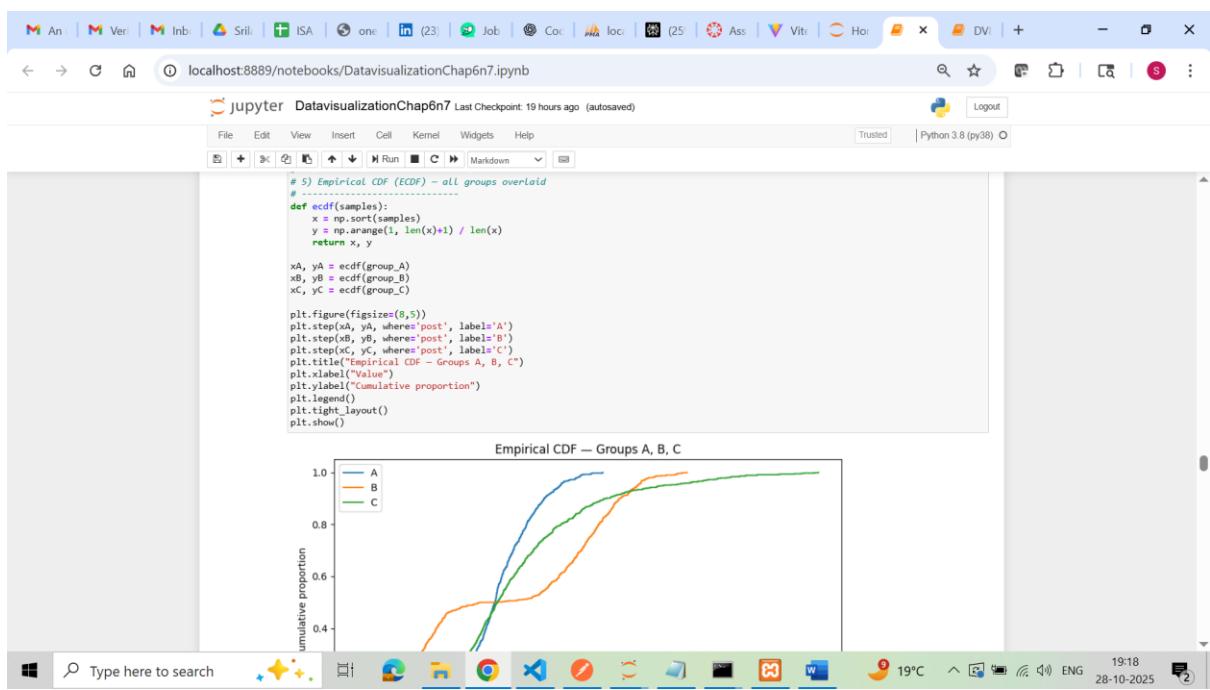
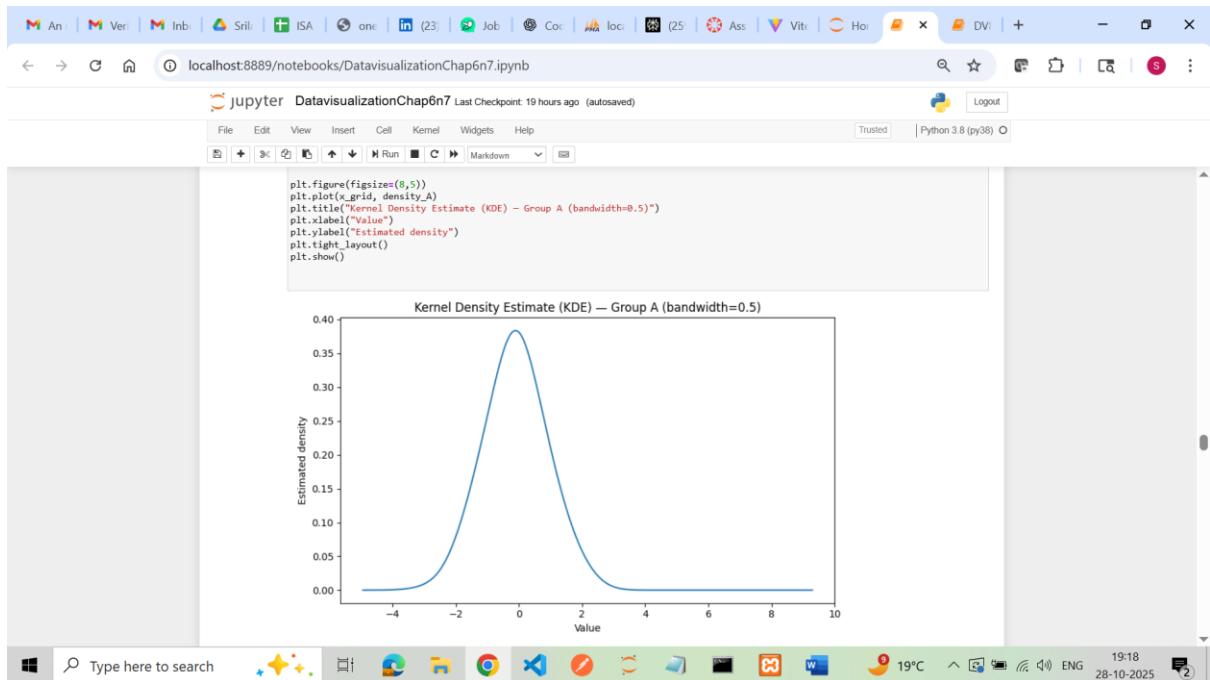


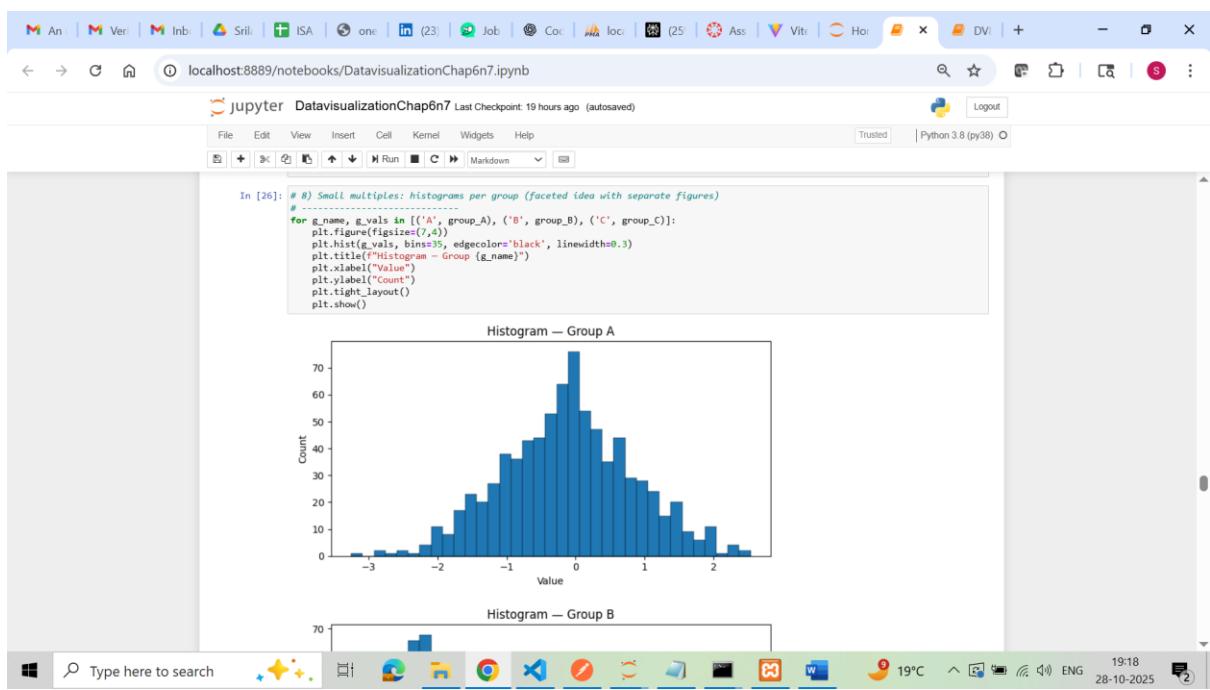
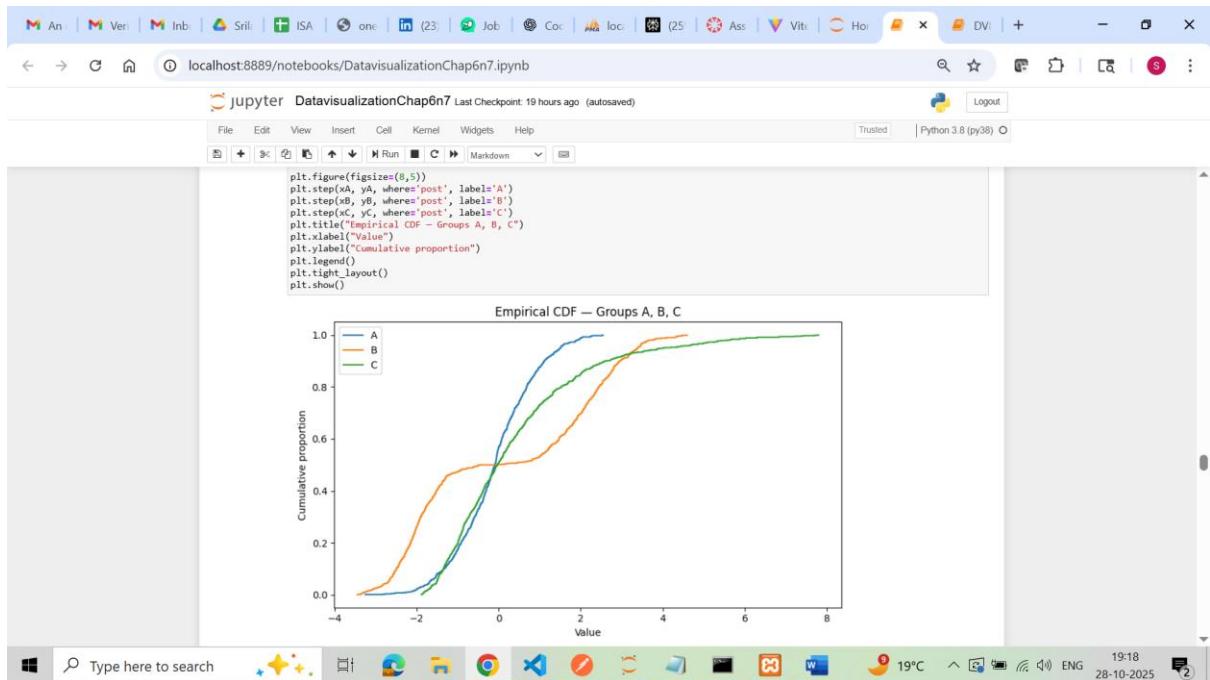
CHAPTER 7

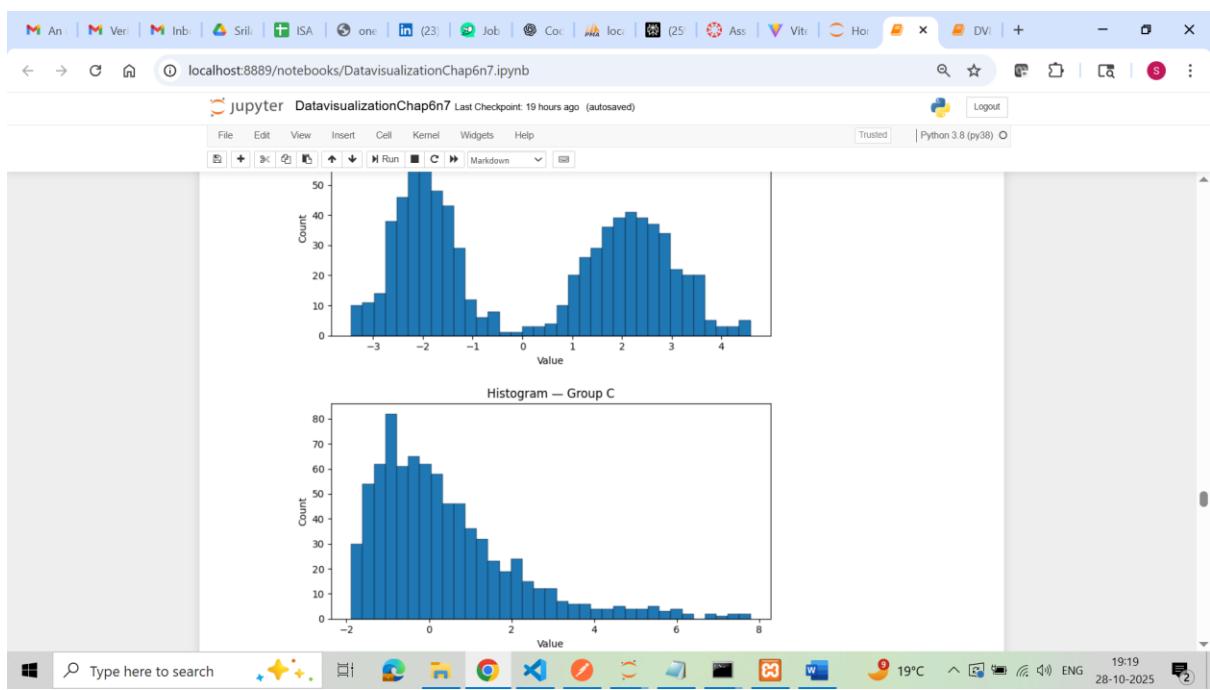
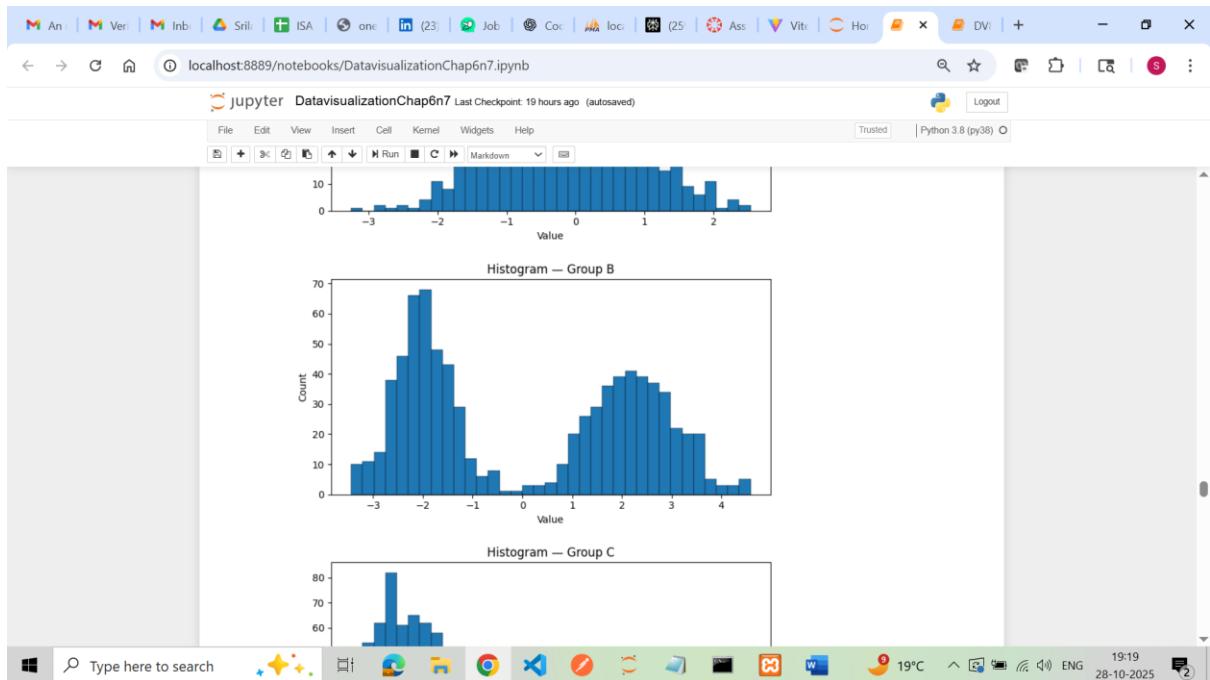












localhost:8889/notebooks/DatavisualizationChap6n7.ipynb

jupyter DatavisualizationChap6n7 Last Checkpoint: 19 hours ago (autosaved)

```

import numpy as np
import matplotlib.pyplot as plt

rng = np.random.default_rng(42)

# -----
# Two-group dataset: session duration (minutes)
# -----
# New users: generally shorter sessions with a long right tail
n_new = 600
session_new = np.concatenate([
    rng.normal(8, 3, size=500),      # main mass around 8 min
    rng.exponential(6, size=100)     # right tail
])
session_new = session_new[session_new > 0]

# Returning users: longer/more variable sessions
n_ret = 1000
session_ret = np.concatenate([
    rng.normal(16, 5, size=800),
    rng.normal(28, 6, size=200)
])
session_ret = session_ret[session_ret > 0]

# Common bins so comparisons are fair
bins = np.linspace(0, max(session_new.max(), session_ret.max()), 30)

# A) BAD: Stacked histogram (easily misread)
# -----
plt.figure(figsize=(9,5))
plt.hist([session_new, session_ret], bins=bins, stacked=True, label=['New','Returning'])
plt.xlabel("Session duration (min)")
plt.ylabel("Count")
plt.title("BAD: Stacked histogram of session durations (New vs Returning)")
plt.legend()
plt.tight_layout()
plt.show()

```

Type here to search 19°C ENG 28-10-2025

localhost:8889/notebooks/DatavisualizationChap6n7.ipynb

jupyter DatavisualizationChap6n7 Last Checkpoint: 19 hours ago (autosaved)

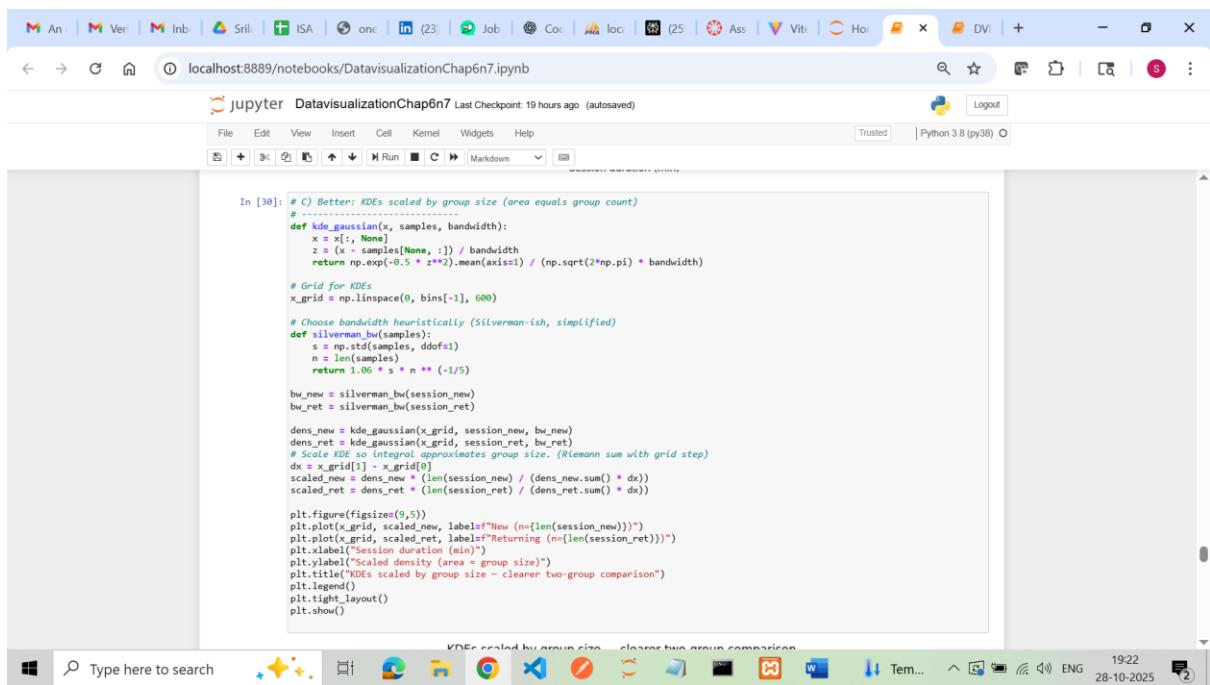
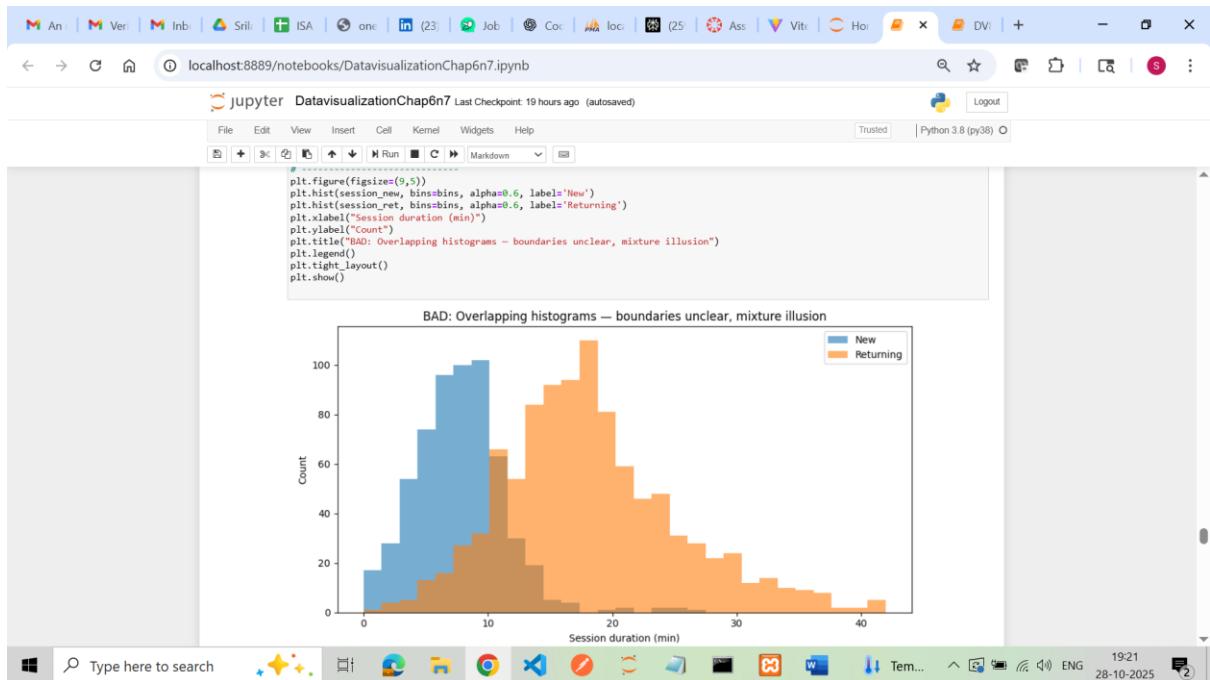
```

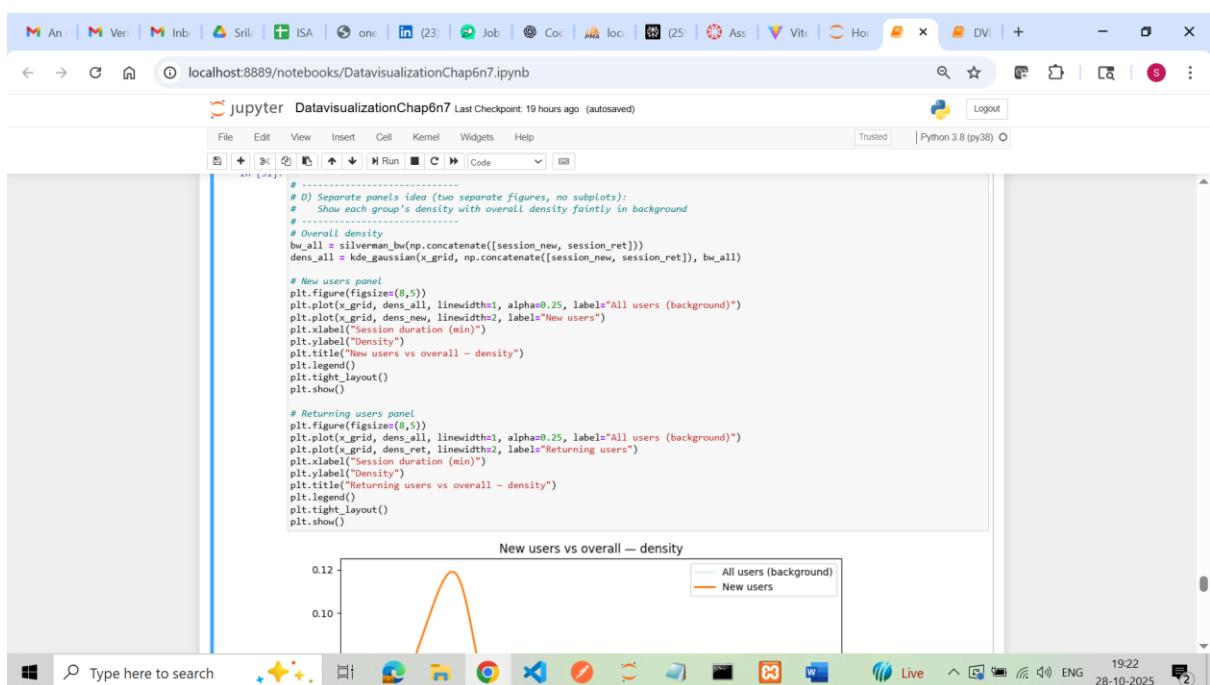
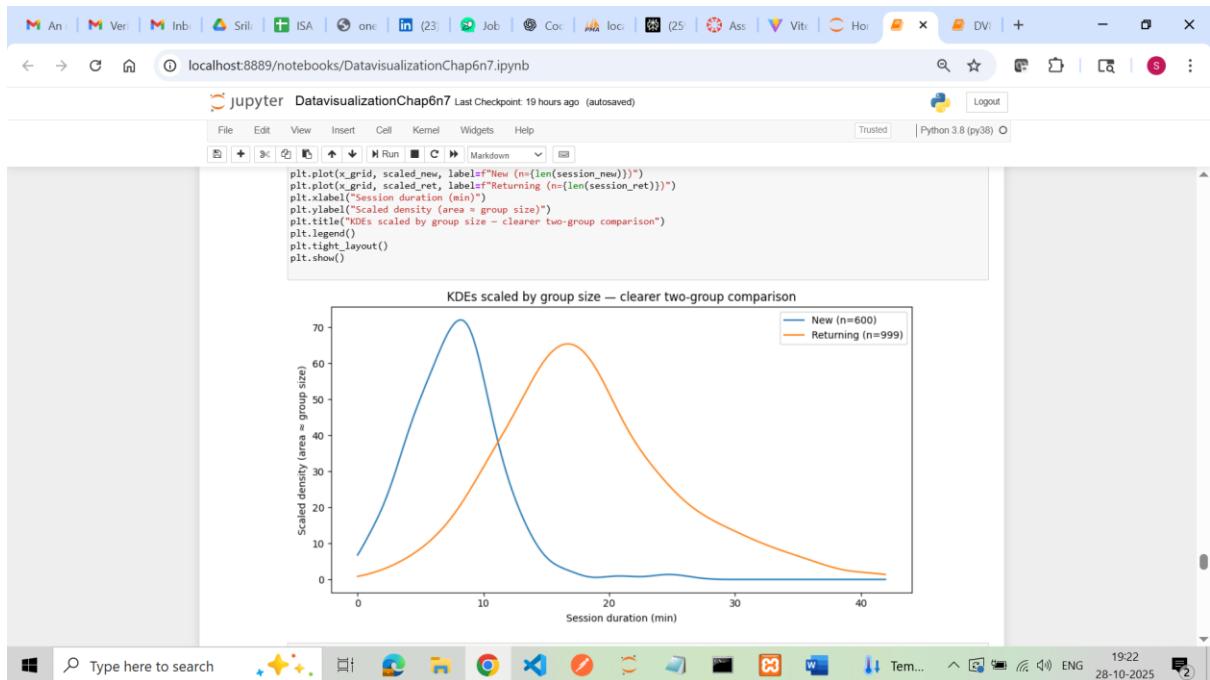
plt.figure(figsize=(9,5))
plt.hist([session_new, session_ret], bins=bins, stacked=True, label=['New','Returning'])
plt.xlabel("Session duration (min)")
plt.ylabel("Count")
plt.title("BAD: Stacked histogram of session durations (New vs Returning)")
plt.legend()
plt.tight_layout()
plt.show()

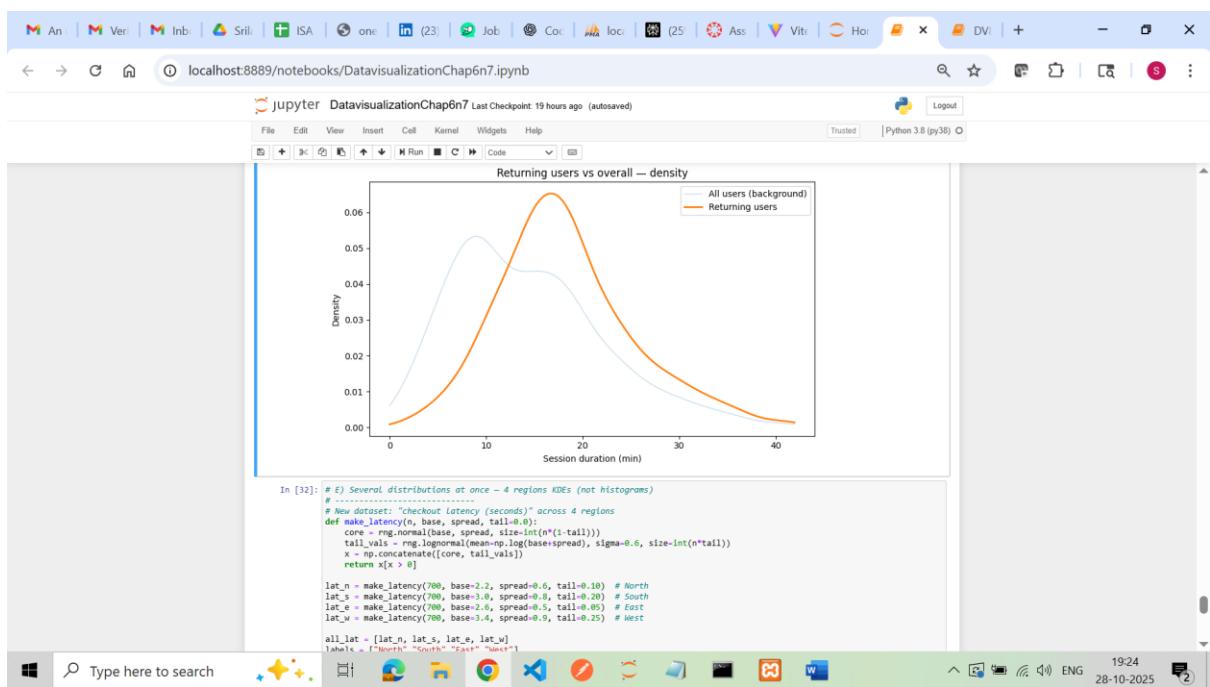
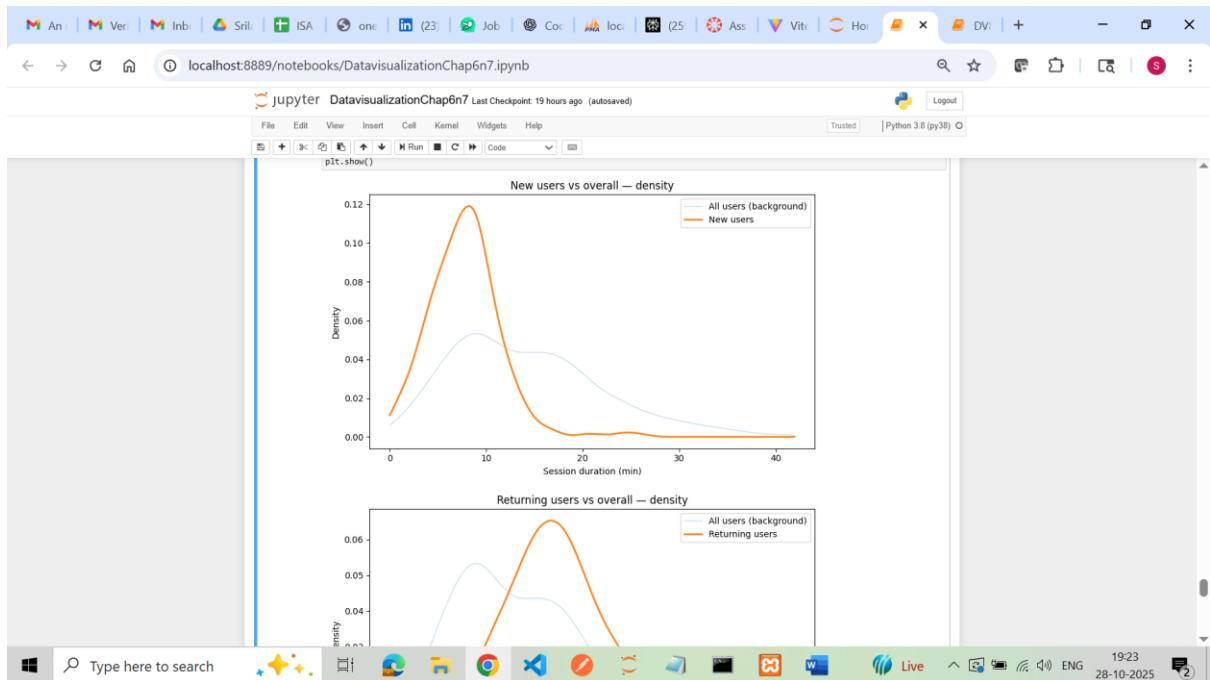
```

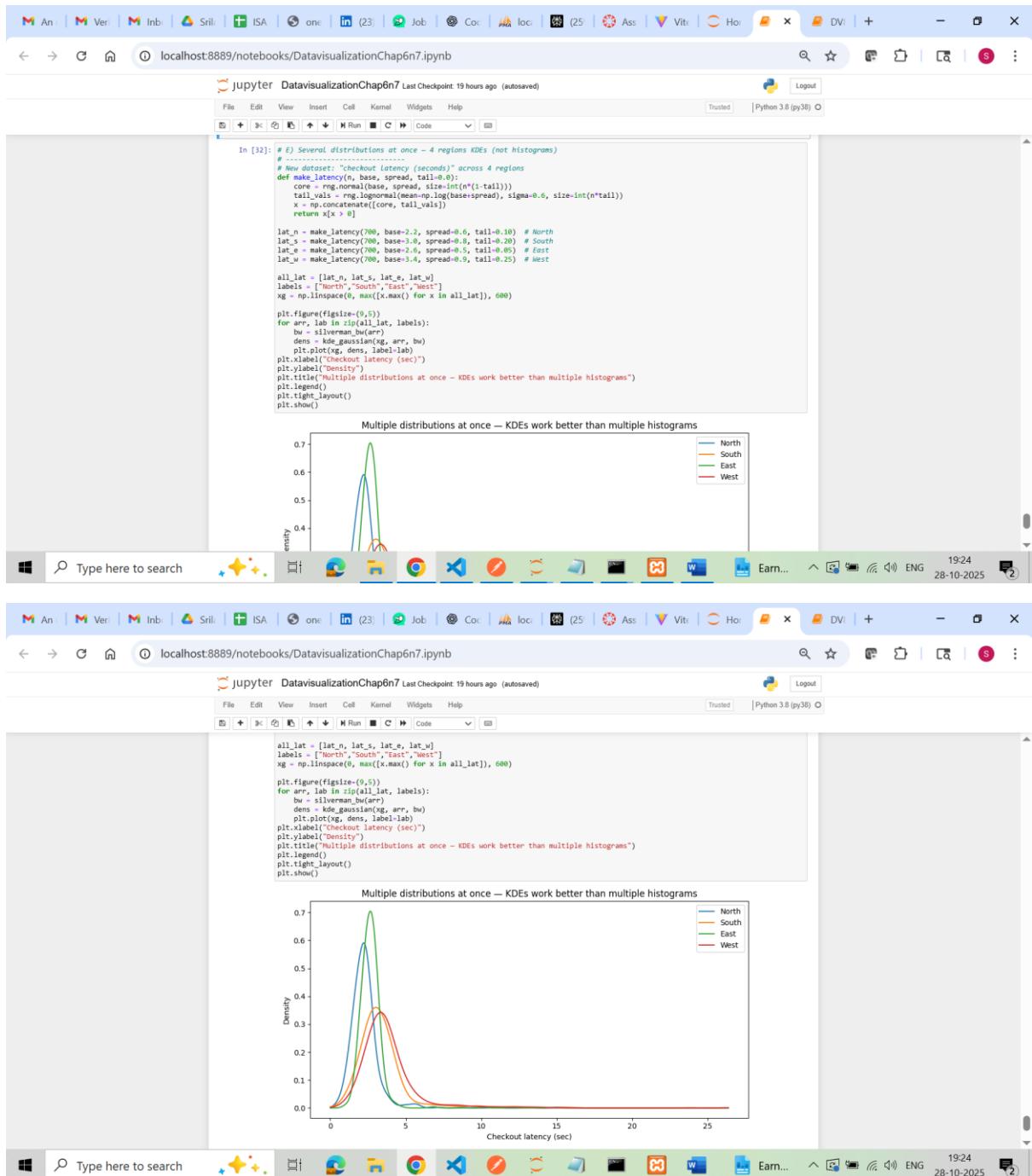
In [28]: # B) Overlapping histograms (alpha) - Looks Like 3 groups

Type here to search 19:20 ENG 28-10-2025









CHAPTER 8

localhost:8889/notebooks/DV8.ipynb

jupyter DV8 Last Checkpoint: 4 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38) Logout

CHAPTER 8

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from scipy import stats

rng = np.random.default_rng(42)

# -----
# Dataset setup
# -----
scores = rng.normal(75, 10, 50) # Exam scores ~ Normal
scores = np.clip(scores, 0, 100)

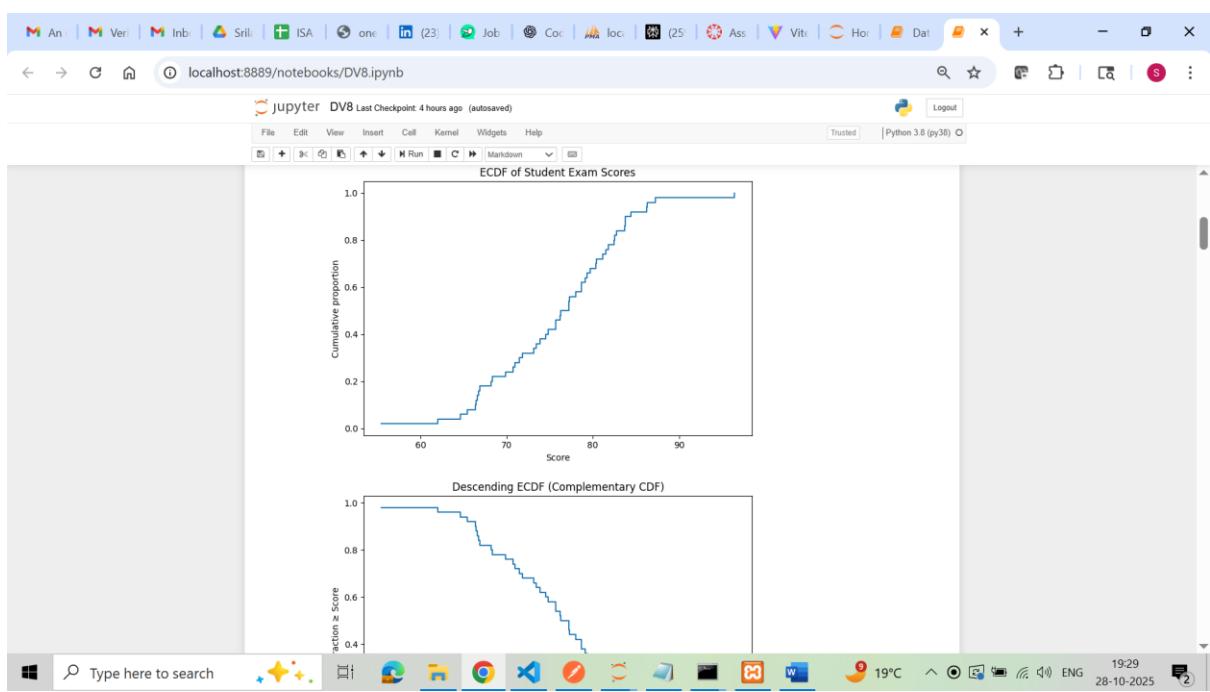
population = np.exp(rng.normal(10, 1.2, 3000)) # Right-skewed (log-normal)
word_counts = np.random.zipf(a=2.0, size=5000) # Very heavy tail

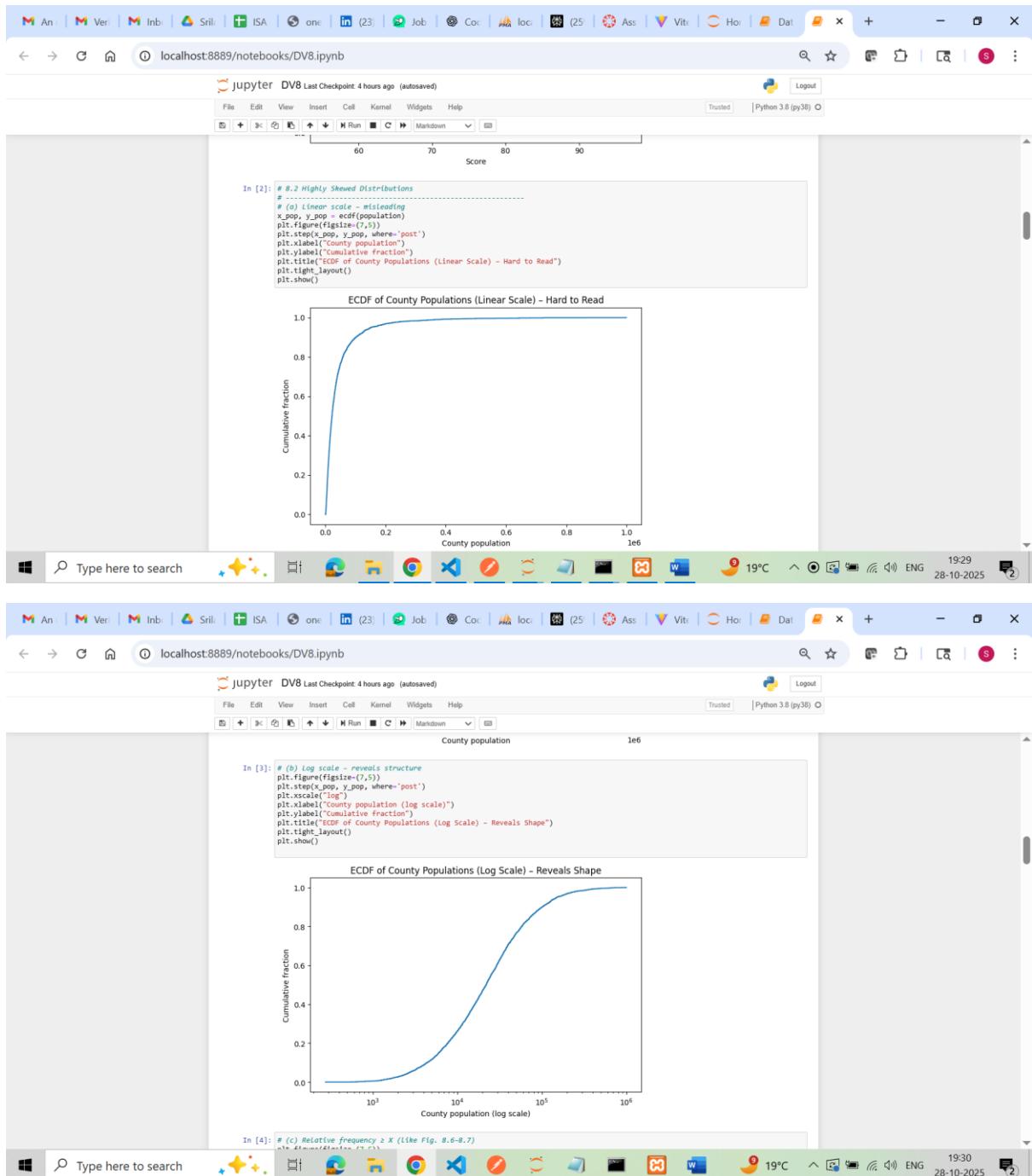
# -----
# 8.1 Empirical Cumulative Distribution Function (ECDF)
# -----
def ecdf(data):
    x = np.sort(data)
    y = np.arange(1, len(x)+1) / len(x)
    return x, y

x_scores, y_scores = ecdf(scores)

plt.figure(figsize=(7,5))
plt.step(x_scores, y_scores, where='post')
plt.xlabel("Score")
plt.ylabel("Cumulative proportion")
plt.title("ECDF of Student Exam Scores")
plt.tight_layout()
plt.show()

# Descending ECDF
plt.figure(figsize=(7,5))
plt.step(x_scores, 1 - y_scores, where='post')
plt.xlabel("Score")
plt.ylabel("Fraction > Score")
plt.title("Descending ECDF (Complementary CDF)")
plt.tight_layout()
plt.show()
```





localhost:8889/notebooks/DV8.ipynb

jupyter DV8 Last Checkpoint: 4 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38) Logout

In [4]: # (c) Relative frequency $\geq X$ (Like Fig. 8.6-8.7)
 $\text{plt.figure(figsize=(7,5)}$
 $\text{counts} = \text{np.sort(population)}$
 $\text{rel_freq} = 1 - \text{np.arange(1, len(counts)+1)} / \text{len(counts)}$
 $\text{plt.loglog(counts, rel_freq)}$
 $\text{plt.xlabel("County population")}$
 $\text{plt.ylabel("Fraction } \geq X")$
 $\text{plt.title("Relative Frequency } \geq X \text{ (Power-Law-like Tail)")}$
 $\text{plt.tight_layout()}$
 plt.show()

County population (log scale)

Relative Frequency $\geq X$ (Power-Law-like Tail)

In [5]: # 8.3 Quantile-Quantile (Q-Q) Plots

Type here to search 19°C ENG 28-10-2025

jupyter DV8 Last Checkpoint: 4 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38) Logout

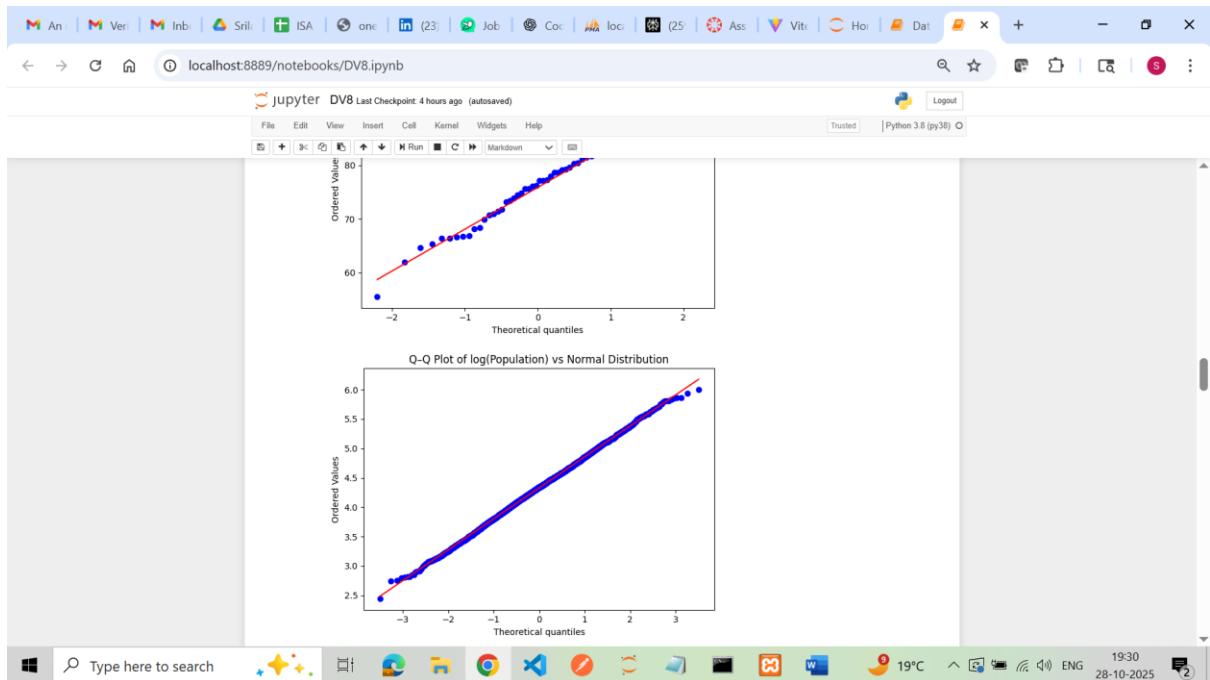
In [5]: # 8.3 Quantile-Quantile (Q-Q) Plots

(a) Q-Q plot of exam scores vs Normal distribution
stats.probplot(scores, dist="norm", plot=plt)
plt.title("Q-Q Plot of Exam Scores vs Normal Distribution")
plt.tight_layout()
plt.show()

(b) Q-Q plot of log(population) vs Normal
log_pop = np.log10(population)
stats.probplot(log_pop, dist="norm", plot=plt)
plt.title("Q-Q Plot of log(Population) vs Normal Distribution")
plt.tight_layout()
plt.show()

Q-Q Plot of Exam Scores vs Normal Distribution

Type here to search 19°C ENG 28-10-2025



CHAPTER 9

