```python
1  # airline_delay_analysis.py
2
3  import pandas as pd
4  import numpy as np
5  import seaborn as sns
6  import matplotlib.pyplot as plt
7
8  # 1. Simulate realistic flight delay data (since
   actual file not available)
9  np.random.seed(42)
10 n = 2000
11 airports = ['DEL', 'BOM', 'BLR', 'MAA', 'CCU']
12 months = np.random.choice(range(1, 13), n)
13 seasons = pd.cut(months, bins=[0, 3, 6, 9, 12],
   labels=['Winter', 'Summer', 'Monsoon', 'Post-Monsoon'
   ])
14
15 data = pd.DataFrame({
16     'FlightNumber': ['AI' + str(i).zfill(4) for i in
   range(n)],
17     'Airport': np.random.choice(airports, n),
18     'ScheduledDeparture': pd.date_range(start='2024-
   01-01', periods=n, freq='H'),
19     'ActualDeparture': pd.date_range(start='2024-01-
   01', periods=n, freq='H') + pd.to_timedelta(np.random
   .randint(0, 180, n), unit='m'),
20     'DelayReason': np.random.choice(['Weather', '
   Carrier', 'ATC', 'Security', 'None'], n, p=[0.25, 0.
   35, 0.2, 0.1, 0.1]),
21     'RouteLength_km': np.random.normal(1200, 300, n).
   astype(int),
22     'FlightFreq_perDay': np.random.choice([1, 2, 3, 4
   ], n),
23     'Month': months,
24     'Season': seasons
25 })
26
27 # 2. Feature Engineering
28 data['DelayMinutes'] = (data['ActualDeparture'] -
   data['ScheduledDeparture']).dt.total_seconds() / 60
29 data = data[data['DelayReason'] != 'None']  # remove
```

```python
29  non-delayed for focused analysis
30  data['DelayCategory'] = pd.cut(data['DelayMinutes'],
    bins=[0, 30, 90, np.inf], labels=['Short', 'Medium',
    'Long'])
31
32  # 3. Exploratory Data Analysis
33  print("\n--- Average Delay by Airport ---")
34  print(data.groupby('Airport')['DelayMinutes'].mean())
35
36  print("\n--- Delay Reason Distribution ---")
37  print(data['DelayReason'].value_counts(normalize=True
    ) * 100)
38
39  # Bar plot of delay reasons
40  plt.figure(figsize=(6, 4))
41  sns.countplot(data=data, x='DelayReason', order=data[
    'DelayReason'].value_counts().index, palette='Set2')
42  plt.title("Delay Reason Distribution")
43  plt.ylabel("Number of Flights")
44  plt.xticks(rotation=30)
45  plt.tight_layout()
46  plt.show()
47
48  # Line plot of monthly delays
49  monthly_delay = data.groupby('Month')['DelayMinutes'
    ].mean()
50  plt.figure(figsize=(7, 4))
51  monthly_delay.plot(marker='o')
52  plt.title("Average Monthly Delay (All Airports)")
53  plt.xlabel("Month")
54  plt.ylabel("Avg Delay (min)")
55  plt.grid()
56  plt.tight_layout()
57  plt.show()
58
59  # Heatmap of delay by airport & month
60  pivot = data.pivot_table(index='Airport', columns='
    Month', values='DelayMinutes', aggfunc='mean')
61  plt.figure(figsize=(8, 5))
62  sns.heatmap(pivot, annot=True, fmt=".1f", cmap='
    coolwarm')
```

```python
63 plt.title("Avg Delay by Airport and Month")
64 plt.tight_layout()
65 plt.show()
66
67 # 4. Root Cause by Season
68 seasonal_cause = pd.crosstab(data['Season'], data['
   DelayReason'], normalize='index') * 100
69 seasonal_cause.plot(kind='bar', stacked=True,
   colormap='tab20', figsize=(8, 4))
70 plt.title("Root Cause Attribution by Season (%)")
71 plt.ylabel("Percentage of Delays")
72 plt.tight_layout()
73 plt.show()
74
75 # 5. Correlation Analysis
76 corr_data = data[['DelayMinutes', 'RouteLength_km',
   'FlightFreq_perDay']]
77 corr = corr_data.corr()
78 sns.heatmap(corr, annot=True, cmap='Blues')
79 plt.title("Correlation Matrix")
80 plt.tight_layout()
81 plt.show()
82
83 # 6. Top 5 Most Delayed Flights
84 top5 = data.sort_values(by='DelayMinutes', ascending
   =False).head(5)
85 print("\n--- Top 5 Most Delayed Flights ---")
86 print(top5[['FlightNumber', 'Airport', 'DelayMinutes
   ', 'DelayReason', 'ScheduledDeparture']])
87
```