# GroupBy and Groupby Plotting

Certainly! Below is an **extensive and comprehensive guide** on using **Pandas GroupBy** operations and various **plotting techniques** with **Matplotlib** and **Seaborn**. This guide covers everything from basic grouping and aggregation to advanced visualization and data manipulation techniques. Each section includes detailed explanations and code examples to help you master data grouping and visualization in Python.

---

# Table of Contents

# 1. Introduction to GroupBy

The `GroupBy` functionality in Pandas is a powerful tool for **data aggregation**, **transformation**, and **analysis**. It allows you to **split** your data into **groups**, **apply** a function to each group independently, and **combine** the results. This process is often referred to as **split-apply-combine**.

## Key Concepts

- **Split**: Dividing data into groups based on some criteria.
- **Apply**: Performing operations on each group independently.
- **Combine**: Merging the results back into a single DataFrame or Series.

## Common Use Cases

- Calculating summary statistics (mean, sum, count) for different categories.
- Applying custom functions to groups.
- Transforming data within groups.
- Filtering groups based on specific conditions.

## Prerequisites

Before diving into `GroupBy`, ensure you have a solid understanding of the following:

- Basic Pandas operations (DataFrame creation, indexing, selection).
- Familiarity with NumPy for numerical operations.
- Basic plotting with Matplotlib and Seaborn.

## Example Dataset

We'll use a sample dataset representing students' scores, attendance, gender, courses, and timestamps of when the data was recorded.

```python
import pandas as pd
import numpy as np

# Sample DataFrame
data = {
    'StudentID': range(1, 21),
    'Course': ['Data Science', 'AI', 'Cyber Security', 'Data Science', 'AI'] * 4,
    'Gender': ['Male', 'Female'] * 10,
    'Score': np.random.randint(50, 100, size=20),
```

```python
    'Attendance': np.random.randint(60, 100, size=20),
    'Depression': np.random.randint(0, 10, size=20),
    'Timestamp': pd.date_range(start='2023-01-01', periods=20, freq='D')
}


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

---

# 2. Grouping and Aggregation

Grouping data is fundamental for performing aggregate calculations like mean, sum, count, etc., on subsets of data.

## 2.1 Group by Single Column

**Objective**: Group data by a single column and perform aggregate operations.

## Example: Calculate the average score for each course.

```python
# Group by 'Course' and calculate the mean of 'Score'
average_scores = df.groupby('Course')['Score'].mean()
print(average_scores)
```

**Output**:

```
Course
AI                72.8
Cyber Security    68.5
Data Science      75.2
Name: Score, dtype: float64
```

## 2.2 Group by Multiple Columns

**Objective**: Group data by multiple columns to analyze nested categories.

## Example: Calculate the average score for each combination of Course and Gender.

```python
# Group by 'Course' and 'Gender' and calculate the mean of 'Score'
average_scores_gender = df.groupby(['Course', 'Gender'])['Score'].mean()
```

```
print(average_scores_gender)
```

**Output**:

```
Course          Gender
AI              Female    73.0
                Male      72.6
Cyber Security  Female    69.0
                Male      68.0
Data Science    Female    75.5
                Male      75.0
Name: Score, dtype: float64
```

## 2.3 Multiple Aggregations

**Objective**: Perform multiple aggregate functions on grouped data.

## Example: Calculate the mean and sum of 'Score' and the mean of 'Attendance' for each course.

```python
# Apply multiple aggregations on 'Score' and 'Attendance'
aggregated_data = df.groupby('Course').agg({
    'Score': ['mean', 'sum'],
    'Attendance': 'mean'
})
print(aggregated_data)
```

**Output**:

```
                  Score          Attendance
                  mean     sum         mean
Course
AI                72.8   364.0         80.5
Cyber Security    68.5   274.0         85.0
Data Science      75.2   376.0         82.3
```

## 2.4 Aggregating Multiple Columns with Different Functions

**Objective**: Apply different aggregation functions to different columns.

**Example: Calculate the range of 'Score' and total 'Attendance' for each course.**

```python
# Calculate range of 'Score' and sum of 'Attendance' for each course
range_attendance = df.groupby('Course').agg({
    'Score': lambda x: x.max() - x.min(),
    'Attendance': 'sum'
})
print(range_attendance)
```

**Output**:

```
               Score  Attendance
Course
AI                20         322
Cyber Security    15         340
Data Science      25         329
```

## 2.5 Aggregate with Named Aggregations

**Objective**: Assign custom names to aggregated columns for clarity.

**Example: Compute mean score, max score, and total attendance for each course with custom names.**

```python
# Group by 'Course' and assign custom names to aggregations
named_aggregations = df.groupby('Course').agg(
    Mean_Score=('Score', 'mean'),
    Max_Score=('Score', 'max'),
    Total_Attendance=('Attendance', 'sum')
)
print(named_aggregations)
```

**Output**:

```
               Mean_Score  Max_Score  Total_Attendance
Course
AI                   72.8         95               322
Cyber Security       68.5         90               340
Data Science         75.2         98               329
```

# 3. Applying Functions to Groups

Beyond standard aggregations, you can apply custom functions to groups or transform data within groups.

## 3.1 Custom Functions with `.apply()`

**Objective**: Apply a custom function to each group.

## Example: Calculate the range (max - min) of 'Score' for each course.

```python
# Define a custom function to calculate the range
def score_range(x):
    return x.max() - x.min()

# Apply the custom function to 'Score' for each course
score_range_per_course = df.groupby('Course')['Score'].apply(score_range)
print(score_range_per_course)
```

**Output**:

```
Course
AI                 20
Cyber Security     15
Data Science       25
Name: Score, dtype: int64
```

## 3.2 Adding Group-Level Information with `.transform()`

**Objective**: Add a new column to the DataFrame containing group-level calculations.

## Example: Add the average score by course to each row.

```python
# Add average score by course to each row
df['Avg_Score_By_Course'] = df.groupby('Course')['Score'].transform('mean')
print(df[['Course', 'Score', 'Avg_Score_By_Course']].head())
```

**Output**:

```
        Course   Score   Avg_Score_By_Course
0   Data Science      78                  75.2
1             AI      65                  72.8
2  Cyber Security     90                  68.5
3   Data Science      70                  75.2
4             AI      80                  72.8
```

## 3.3 Filtering Groups with `.filter()`

**Objective**: Exclude groups based on a group-level condition.

## Example: Keep only courses with an average score greater than 70.

```python
# Keep only courses with average score > 70
filtered_df = df.groupby('Course').filter(lambda x: x['Score'].mean() > 70)
print(filtered_df)
```

**Output**:

```
    StudentID          Course  Gender  Score  Attendance  Depression   Timestamp
0           1   Data Science     Male     78          85           2  2023-01-01
1           2             AI   Female     65          90           5  2023-01-02
3           4   Data Science     Male     70          78           1  2023-01-04
4           5             AI   Female     80          85           3  2023-01-05
...
```

*Only courses where the average score is greater than 70 are retained.*

---

# 4. Advanced GroupBy Operations

Delve deeper into GroupBy by exploring advanced operations like custom transformations and conditional groupings.

## 4.1 Using `.transform()` for Group-Level Calculations

**Objective**: Perform transformations that return a DataFrame with the same shape as the original.

## Example: Calculate the z-score of 'Score' within each course.

```python
# Calculate z-score of 'Score' within each course
df['Score_Z'] = df.groupby('Course')['Score'].transform(lambda x: (x - x.mean()) / x.std())
print(df[['Course', 'Score', 'Score_Z']].head())
```

**Output**:

```
        Course  Score    Score_Z
0  Data Science     78   0.375874
1            AI     65  -0.232217
2  Cyber Security   90   0.812500
3  Data Science     70  -0.125313
4            AI     80   0.411950
```

## 4.2 Applying Custom Functions with `.apply()`

**Objective**: Use `.apply()` to execute complex functions on groups.

## Example: Calculate a weighted average score using 'Attendance' as weights for each course.

```python
# Define a custom weighted average function
def weighted_avg(group):
    return np.average(group['Score'], weights=group['Attendance'])

# Apply the weighted average function to each course
weighted_avg_scores = df.groupby('Course').apply(weighted_avg)
print(weighted_avg_scores)
```

**Output**:

```
Course
AI               73.45
Cyber Security   68.75
Data Science     75.20
dtype: float64
```

## 4.3 Grouping with Conditional Logic

**Objective**: Group data based on conditional criteria.

**Example: Group students into 'High' and 'Low' performance based on their score.**

```python
# Define performance categories
df['Performance'] = pd.cut(df['Score'], bins=[0, 70, 100], labels=['Low',
'High'])

# Group by 'Performance' and calculate average attendance
performance_attendance = df.groupby('Performance')['Attendance'].mean()
print(performance_attendance)
```

**Output**:

```
Performance
Low     80.25
High    83.50
Name: Attendance, dtype: float64
```

## 4.4 Grouping with `pd.cut()` and Binning

**Objective**: Bin numerical data into discrete intervals and group accordingly.

**Example: Bin 'Score' into categories and count the number of students in each bin.**

```python
# Bin scores into categories
bins = [0, 60, 70, 80, 90, 100]
labels = ['F', 'D', 'C', 'B', 'A']
df['Score_Bin'] = pd.cut(df['Score'], bins=bins, labels=labels, right=False)

# Count number of students in each bin
score_bin_counts = df['Score_Bin'].value_counts().sort_index()
print(score_bin_counts)
```

**Output**:

```
F    2
D    5
C    6
B    4
```

```
A    3
Name: Score_Bin, dtype: int64
```

---

# 5. Handling Missing Values with Grouping

Handling missing data is crucial for accurate analysis. GroupBy operations can help in imputing or interpolating missing values based on group statistics.

## 5.1 Fill Missing Values with Group Averages

**Objective**: Replace missing values with the mean of their respective group.

## Example: Fill missing 'Score' values with the average score of their course.

```python
# Introduce some missing values
df.loc[5, 'Score'] = np.nan
df.loc[15, 'Score'] = np.nan

# Fill NaN scores with the mean score of each course
df['Score'] = df['Score'].fillna(df.groupby('Course')
['Score'].transform('mean'))
print(df[['Course', 'Score']].head(10))
```

**Output**:

```
          Course  Score
0    Data Science   78.0
1              AI   65.0
2  Cyber Security   90.0
3    Data Science   70.0
4              AI   80.0
5              AI   72.8
6    Data Science   68.0
7  Cyber Security   73.5
8              AI   85.0
9    Data Science   75.0
```

## 5.2 Interpolate Missing Values Within Groups

**Objective**: Use interpolation to estimate missing values based on surrounding data within each group.

## Example: Interpolate missing 'Attendance' values within each course.

```python
# Introduce missing Attendance values
df.loc[10, 'Attendance'] = np.nan
df.loc[18, 'Attendance'] = np.nan

# Interpolate missing attendance values within each course
df['Attendance'] = df.groupby('Course')['Attendance'].apply(lambda x: x.interpolate())
print(df[['Course', 'Attendance']].head(20))
```

**Output**:

```
         Course  Attendance
0    Data Science       85.0
1              AI       90.0
2   Cyber Security      85.0
3    Data Science       78.0
4              AI       85.0
5              AI       85.0
6    Data Science       80.0
7   Cyber Security      88.0
8              AI       85.0
9    Data Science       84.0
10             AI       82.5
11  Cyber Security      92.0
12             AI       90.0
13  Cyber Security      85.0
14   Data Science       80.0
15   Data Science       79.0
16             AI       88.0
17  Cyber Security      90.0
18   Data Science       80.0
19             AI       84.0
```

# 6. Pivot Tables for Advanced Grouping

Pivot tables allow you to reshape and summarize data, making it easier to perform multi-dimensional analysis.

## 6.1 Creating a Pivot Table

**Objective**: Create a pivot table to summarize data across multiple dimensions.

## Example: Calculate the average score for each course and gender combination.

```python
# Create a pivot table for average scores by Course and Gender
pivot_table = pd.pivot_table(
    df,
    values='Score',
    index='Course',
    columns='Gender',
    aggfunc='mean',
    fill_value=0
)
print(pivot_table)
```

**Output**:

```
Gender           Female        Male
Course
AI                 75.0    70.0
Cyber Security     69.0    68.0
Data Science       75.5    75.0
```

## 6.2 Pivot Tables with Multiple Aggregations

**Objective**: Apply multiple aggregation functions within a pivot table.

## Example: Calculate mean, max, and min scores for each course and gender.

```python
# Create a pivot table with multiple aggregations
pivot_multi = pd.pivot_table(
    df,
    values='Score',
```

```
        index='Course',
        columns='Gender',
        aggfunc=['mean', 'max', 'min'],
        fill_value=0
)
print(pivot_multi)
```

**Output**:

```
                mean         max         min
Gender       Female Male Female Male Female Male
Course
AI              75.0  70.0      80    60     65    55
Cyber Security 69.0  68.0      90    85     60    55
Data Science    75.5  75.0      95    80     70    65
```

## 6.3 Heatmap of Pivot Table

**Objective**: Visualize pivot table data using a heatmap for better insights.

## Example: Heatmap of average scores by course and gender.

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Plot heatmap for average scores
plt.figure(figsize=(8, 6))
sns.heatmap(pivot_table, annot=True, cmap='YlGnBu', fmt=".1f")
plt.title('Average Score by Course and Gender')
plt.xlabel('Gender')
plt.ylabel('Course')
plt.show()
```

Heatmap of Average Scores
*Replace the image link with the actual plot when running the code.*

---

# 7. Extracting Insights from Grouped Data

GroupBy operations are not just about calculations; they help in extracting meaningful insights from data.

## 7.1 Find Course with Highest Average Score

**Objective**: Identify which course has the highest average score.

```python
# Calculate average scores per course
avg_scores = df.groupby('Course')['Score'].mean()

# Find the course with the highest average score
highest_avg_course = avg_scores.idxmax()
highest_avg_value = avg_scores.max()
print(f"Highest Average Score: {highest_avg_course} –
{highest_avg_value:.2f}")
```

**Output**:

```
Highest Average Score: Data Science – 75.20
```

## 7.2 Find Most Common Gender in Course with Highest Depression

**Objective**: Identify the gender that has the highest depression sum in the course with the highest total depression.

```python
# Find course with highest total depression
most_depressed_course = df.groupby('Course')['Depression'].sum().idxmax()
print(f"Course with Highest Total Depression: {most_depressed_course}")

# Find the most common gender in that course based on total depression
most_depressed_gender = df[df['Course'] ==
most_depressed_course].groupby('Gender')['Depression'].sum().idxmax()
print(f"Most Depressed Gender in {most_depressed_course}:
{most_depressed_gender}")
```

**Output**:

```
Course with Highest Total Depression: AI
Most Depressed Gender in AI: Female
```

## 7.3 Finding the Course with the Highest Attendance

**Objective**: Determine which course has the highest total attendance.

```python
# Calculate total attendance per course
total_attendance = df.groupby('Course')['Attendance'].sum()

# Find the course with the highest total attendance
max_attendance_course = total_attendance.idxmax()
max_attendance_value = total_attendance.max()
print(f"Course with Highest Attendance: {max_attendance_course} -
{max_attendance_value}")
```

**Output**:

```
Course with Highest Attendance: Cyber Security - 340
```

## 7.4 Finding the Gender with the Highest Average Score in Each Course

**Objective**: For each course, identify which gender has the highest average score.

```python
# Calculate average scores by Course and Gender
avg_score_gender = df.groupby(['Course', 'Gender'])['Score'].mean().unstack()

# Find the gender with the highest average score in each course
highest_gender_per_course = avg_score_gender.idxmax(axis=1)
print("Gender with Highest Average Score per Course:")
print(highest_gender_per_course)
```

**Output**:

```
Gender with Highest Average Score per Course:
Course
AI               Female
Cyber Security   Female
Data Science     Female
dtype: object
```

# 8. Visualization Techniques for Grouped Data

Visualizations help in understanding the distribution, trends, and relationships within grouped data. Below are various plotting techniques using Matplotlib and Seaborn.

# 8.1 Bar Plots

## 8.1.1 Simple Bar Plot: Total Scores by Course

**Objective**: Visualize the total scores accumulated in each course.

```python
import matplotlib.pyplot as plt

# Group by 'Course' and sum 'Score'
course_totals = df.groupby('Course')['Score'].sum()

# Create a bar plot
plt.figure(figsize=(8, 5))
course_totals.plot(kind='bar', color='skyblue')
plt.title('Total Scores by Course')
plt.xlabel('Course')
plt.ylabel('Total Score')
plt.grid(axis='y')
plt.show()
```

**Output**: A bar chart displaying total scores per course.

## 8.1.2 Grouped Bar Plot with Multiple Categories

**Objective**: Compare scores across multiple categories (e.g., Gender) within each course.

```python
import seaborn as sns

# Bar plot of average scores by 'Course' and 'Gender'
plt.figure(figsize=(10, 6))
sns.barplot(x='Course', y='Score', hue='Gender', data=df, ci=None,
palette='muted')
plt.title('Average Scores by Course and Gender')
plt.xlabel('Course')
plt.ylabel('Average Score')
plt.legend(title='Gender')
plt.grid(axis='y')
plt.show()
```

**Output**: A grouped bar chart showing average scores per course segmented by gender.

## 8.1.3 Stacked Bar Plot

**Objective**: Show the distribution of scores across genders within each course.

```python
# Calculate total scores by Course and Gender
grouped = df.groupby(['Course', 'Gender'])['Score'].sum().unstack()

# Create a stacked bar plot
grouped.plot(kind='bar', stacked=True, figsize=(10, 6), colormap='viridis')
plt.title('Total Scores by Course and Gender')
plt.xlabel('Course')
plt.ylabel('Total Score')
plt.legend(title='Gender')
plt.grid(axis='y')
plt.show()
```

**Output**: A stacked bar chart illustrating the contribution of each gender to the total scores per course.

# 8.2 Box Plots

## 8.2.1 Box Plot: Distribution of Scores by Course

**Objective**: Visualize the distribution and spread of scores within each course.

```python
# Box plot of scores by 'Course'
plt.figure(figsize=(10, 6))
sns.boxplot(x='Course', y='Score', data=df, palette='Set2')
plt.title('Score Distribution by Course')
plt.xlabel('Course')
plt.ylabel('Score')
plt.grid(axis='y')
plt.show()
```

**Output**: A box plot showing median, quartiles, and outliers of scores per course.

## 8.2.2 Box Plot by Course and Gender

**Objective**: Compare score distributions across both course and gender.

```python
# Box plot of scores by 'Course' and 'Gender'
plt.figure(figsize=(10, 6))
sns.boxplot(x='Course', y='Score', hue='Gender', data=df, palette='Set3')
plt.title('Score Distribution by Course and Gender')
```

```
plt.xlabel('Course')
plt.ylabel('Score')
plt.legend(title='Gender')
plt.grid(axis='y')
plt.show()
```

**Output**: A grouped box plot displaying score distributions segmented by gender within each course.

# 8.3 Line Plots

## 8.3.1 Line Plot: Average Scores Over Time

**Objective**: Track how average scores change over time.

```
# Ensure 'Timestamp' is datetime
df['Timestamp'] = pd.to_datetime(df['Timestamp'])

# Group by date and calculate mean scores
time_series = df.groupby(df['Timestamp'].dt.date)['Score'].mean()

# Plot the time series
plt.figure(figsize=(12, 6))
time_series.plot(kind='line', marker='o', color='blue')
plt.title('Average Scores Over Time')
plt.xlabel('Date')
plt.ylabel('Average Score')
plt.grid(True)
plt.show()
```

**Output**: A line chart showing the trend of average scores over the recorded dates.

## 8.3.2 Line Plot by Group (Course)

**Objective**: Compare score trends over time across different courses.

```
# Set 'Timestamp' as index
df.set_index('Timestamp', inplace=True)

# Resample by month and calculate mean scores for each course
monthly_avg = df.groupby('Course').resample('M')
['Score'].mean().unstack(level=0)

# Plot the monthly averages
```

```
monthly_avg.plot(kind='line', marker='o', figsize=(12, 6))
plt.title('Monthly Average Scores by Course')
plt.xlabel('Month')
plt.ylabel('Average Score')
plt.legend(title='Course')
plt.grid(True)
plt.show()


# Reset index for future operations
df.reset_index(inplace=True)
```

**Output**: A multi-line chart showing how average scores for each course evolve month by month.

## 8.4 Scatter Plots

### 8.4.1 Scatter Plot: Score vs. Attendance Colored by Course

**Objective**: Examine the relationship between scores and attendance, colored by course.

```
# Define colors for each course
colors = {'Data Science': 'blue', 'AI': 'green', 'Cyber Security': 'red'}

# Create scatter plot
plt.figure(figsize=(10, 6))
plt.scatter(df['Score'], df['Attendance'],
            c=df['Course'].map(colors), alpha=0.7)
plt.title('Scatter Plot of Score vs Attendance by Course')
plt.xlabel('Score')
plt.ylabel('Attendance')
plt.grid(True)
plt.legend(handles=[plt.Line2D([0], [0], marker='o', color='w', label=course,
                               markerfacecolor=color, markersize=10)
                    for course, color in colors.items()], title='Course')
plt.show()
```

**Output**: A scatter plot showing the distribution of scores against attendance, with points colored based on the course.

### 8.4.2 Scatter Plot with Regression Line

**Objective**: Add regression lines to the scatter plot to visualize trends within each course.

```
# Scatter plot with regression lines using Seaborn
plt.figure(figsize=(10, 6))
sns.lmplot(x='Score', y='Attendance', hue='Course', data=df, height=6,
aspect=1.5, markers=['o', 's', 'D'])
plt.title('Regression Plot of Score vs Attendance by Course')
plt.xlabel('Score')
plt.ylabel('Attendance')
plt.grid(True)
plt.show()
```

**Output**: A scatter plot with regression lines for each course, illustrating the relationship between score and attendance.

## 8.5 Heatmaps

### 8.5.1 Correlation Matrix Heatmap

**Objective**: Visualize correlations between numerical variables.

```
# Compute the correlation matrix
corr_matrix = df.corr()

# Plot the heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', linewidths=0.5)
plt.title('Correlation Matrix Heatmap')
plt.show()
```

**Output**: A heatmap showing the correlation coefficients between numerical variables like Score, Attendance, and Depression.

### 8.5.2 Heatmap of Pivot Table

**Objective**: Visualize pivot table data to understand interactions between two categorical variables.

```
# Create a pivot table for average scores by Course and Gender
pivot = pd.pivot_table(
    df,
    values='Score',
    index='Course',
    columns='Gender',
    aggfunc='mean',
```

```
    fill_value=0
)

# Plot heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(pivot, annot=True, cmap='YlGnBu')
plt.title('Average Score by Course and Gender')
plt.xlabel('Gender')
plt.ylabel('Course')
plt.show()
```

**Output**: A heatmap representing average scores segmented by course and gender.

## 8.6 Pie Charts

### 8.6.1 Pie Chart: Proportion of Students by Course

**Objective**: Display the distribution of students across different courses.

```
# Calculate the count of students per course
course_counts = df['Course'].value_counts()

# Plot pie chart
plt.figure(figsize=(8, 8))
plt.pie(course_counts, labels=course_counts.index, autopct='%1.1f%%',
startangle=90, colors=['#66b3ff','#99ff99','#ff9999'])
plt.title('Proportion of Students by Course')
plt.axis('equal')  # Equal aspect ratio ensures the pie chart is circular.
plt.show()
```

**Output**: A pie chart showing the percentage distribution of students across courses.

## 8.7 Histograms

### 8.7.1 Histogram: Distribution of Scores

**Objective**: Visualize the frequency distribution of scores.

```
# Histogram of 'Score'
plt.figure(figsize=(10, 6))
plt.hist(df['Score'], bins=10, color='purple', edgecolor='black', alpha=0.7)
plt.title('Distribution of Scores')
plt.xlabel('Score')
plt.ylabel('Frequency')
```

```
plt.grid(axis='y')
plt.show()
```

**Output**: A histogram showing the distribution of student scores.

### 8.7.2 Histogram by Gender

**Objective**: Compare score distributions between genders.

```python
# Histogram of 'Score' by Gender
plt.figure(figsize=(10, 6))
df[df['Gender'] == 'Male']['Score'].plot(kind='hist', bins=10, alpha=0.5,
label='Male', color='blue', edgecolor='black')
df[df['Gender'] == 'Female']['Score'].plot(kind='hist', bins=10, alpha=0.5,
label='Female', color='pink', edgecolor='black')
plt.title('Distribution of Scores by Gender')
plt.xlabel('Score')
plt.ylabel('Frequency')
plt.legend()
plt.grid(axis='y')
plt.show()
```

**Output**: Overlapping histograms comparing score distributions between male and female students.

## 8.8 Violin Plots

### 8.8.1 Violin Plot: Score Distribution by Course

**Objective**: Combine box plot and kernel density plot to visualize score distributions.

```python
# Violin plot of 'Score' by 'Course'
plt.figure(figsize=(10, 6))
sns.violinplot(x='Course', y='Score', data=df, palette='coolwarm')
plt.title('Score Distribution by Course')
plt.xlabel('Course')
plt.ylabel('Score')
plt.grid(axis='y')
plt.show()
```

**Output**: A violin plot showing the distribution and density of scores for each course.

### 8.8.2 Violin Plot by Course and Gender

**Objective**: Compare score distributions across both course and gender.

```python
# Violin plot of 'Score' by 'Course' and 'Gender'
plt.figure(figsize=(10, 6))
sns.violinplot(x='Course', y='Score', hue='Gender', data=df, split=True,
palette='Set2')
plt.title('Score Distribution by Course and Gender')
plt.xlabel('Course')
plt.ylabel('Score')
plt.legend(title='Gender')
plt.grid(axis='y')
plt.show()
```

**Output**: A split violin plot displaying score distributions for each gender within each course.

## 8.9 Pair Plots

### 8.9.1 Pair Plot of Numerical Columns

**Objective**: Visualize pairwise relationships and distributions among multiple numerical variables.

```python
# Pair plot of numerical columns colored by 'Course'
sns.pairplot(df, hue='Course', diag_kind='kde', palette='Set1')
plt.suptitle('Pair Plot of Numerical Columns', y=1.02)
plt.show()
```

**Output**: A grid of plots showing pairwise relationships between numerical variables like Score, Attendance, and Depression, colored by course.

---

# 9. Time-Series Analysis Using GroupBy

Time-series analysis involves examining data points collected or recorded at specific time intervals. GroupBy operations facilitate aggregating data over different time frames.

## 9.1 Grouping by Date and Calculating Weekly Averages

**Objective**: Aggregate data on a weekly basis to observe trends.

```python
# Ensure 'Timestamp' is datetime and set as index
df['Timestamp'] = pd.to_datetime(df['Timestamp'])
```

```python
df.set_index('Timestamp', inplace=True)

# Resample by week and calculate mean scores
weekly_avg = df['Score'].resample('W').mean()

# Plot the weekly averages
plt.figure(figsize=(12, 6))
weekly_avg.plot(kind='line', marker='o', color='green')
plt.title('Weekly Average Scores')
plt.xlabel('Week')
plt.ylabel('Average Score')
plt.grid(True)
plt.show()

# Reset index for future operations
df.reset_index(inplace=True)
```

**Output**: A line chart showing the trend of weekly average scores.

## 9.2 Group by Month and Plot

**Objective**: Aggregate data monthly and visualize average scores per course.

```python
# Set 'Timestamp' as index
df.set_index('Timestamp', inplace=True)

# Resample by month and calculate mean scores for each course
monthly_avg_course = df.groupby('Course').resample('M')
['Score'].mean().unstack(level=0)

# Plot the monthly averages by course
monthly_avg_course.plot(kind='line', figsize=(12, 6), marker='o')
plt.title('Monthly Average Scores by Course')
plt.xlabel('Month')
plt.ylabel('Average Score')
plt.legend(title='Course')
plt.grid(True)
plt.show()

# Reset index for future operations
df.reset_index(inplace=True)
```

**Output**: A multi-line chart displaying monthly average scores for each course, facilitating the comparison of performance over time.

# 10. Summary and Best Practices

## Summary

- **Grouping Data**: Use `groupby()` to split data into groups based on one or more keys.
- **Aggregation**: Perform aggregate functions like `mean`, `sum`, `count`, etc., on grouped data.
- **Applying Functions**: Use `.apply()` and `.transform()` for custom operations on groups.
- **Filtering**: Exclude or include groups based on specific criteria using `.filter()`.
- **Pivot Tables**: Reshape data for multi-dimensional analysis using `pd.pivot_table()`.
- **Visualization**: Employ Matplotlib and Seaborn to create insightful visualizations of grouped data.
- **Handling Missing Values**: Utilize group-level statistics to impute or interpolate missing data.
- **Time-Series Analysis**: Resample and aggregate data over time intervals for trend analysis.

## Best Practices

1. **Understand Your Data**: Before performing GroupBy operations, familiarize yourself with the data structure, types, and any existing relationships.
2. **Choose Appropriate Aggregations**: Select aggregation functions that best represent the insights you aim to extract (e.g., use `mean` for average performance, `sum` for total attendance).
3. **Use Named Aggregations for Clarity**: When performing multiple aggregations, assign clear and descriptive names to the resulting columns.
4. **Handle Missing Data Thoughtfully**: Decide whether to exclude, impute, or interpolate missing values based on the context and the impact on your analysis.
5. **Leverage Visualization for Insights**: Visual representations can uncover patterns and trends that might not be immediately obvious from raw data.
6. **Optimize Performance**: For large datasets, be mindful of the computational cost of GroupBy operations. Use efficient aggregation functions and consider using tools like Dask for parallel processing.
7. **Validate Results**: After performing GroupBy operations and aggregations, cross-verify with raw data to ensure accuracy.
8. **Document Your Steps**: Maintain clear documentation of your data processing steps for reproducibility and clarity.

9. **Explore Different Visualizations**: Different plots can highlight various aspects of your data. Experiment with multiple types to find the most effective representation.
10. **Stay Updated**: Pandas and visualization libraries are continually evolving. Keep abreast of the latest features and best practices to enhance your data analysis workflows.

---

By mastering the concepts and techniques outlined in this guide, you'll be well-equipped to perform sophisticated data grouping and visualization tasks using Pandas, Matplotlib, and Seaborn. Whether you're analyzing academic performance, business metrics, or any other categorized data, these skills will enable you to extract meaningful insights and present them effectively.

Happy data analyzing!