

How to Easily Create Grouped Boxplots with Matplotlib

Authored by
stats writer

December 5, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Create Grouped Boxplots with Matplotlib*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=105929>

The [Matplotlib](#) library provides robust tools for generating static, animated, and interactive visualizations in Python. While [Matplotlib](#) offers its own dedicated function, the integration with [Seaborn](#) significantly enhances the ease and aesthetic quality of complex plots, such as [boxplots](#) grouped by categorical variables. A grouped [boxplot](#) is a fundamental tool in exploratory [data visualization](#), allowing analysts to quickly compare the distribution, central tendency, and variability of numerical data across distinct categories or groups. Using the powerful `boxplot()` function, often leveraged through the [Seaborn](#) wrapper, we can efficiently visualize these comparative statistics.

Introduction to Grouped Boxplots

Grouped [boxplots](#) (or box-and-whisker plots) are indispensable for understanding the distribution of a numerical variable broken down by one or more categorical variables. Each individual box within the plot summarizes the five-number summary--minimum, first quartile (Q1), median (Q2), third quartile (Q3), and maximum--for a specific group. This visual representation is extremely effective for identifying differences in spread, skewness, and the presence of [outliers](#) between different populations or experimental conditions. By grouping these plots together, we gain immediate comparative insight into how the central tendency and dispersion of the data vary across distinct categories.

When working within the Python ecosystem, achieving high-quality grouped visualizations is streamlined by combining the foundational plotting capabilities of [Matplotlib](#) with the high-level interface offered by [Seaborn](#). [Seaborn](#) is built on top of [Matplotlib](#) and is designed to make statistical graphics appealing and informative with minimal effort. While direct [Matplotlib](#) implementation is possible, the `seaborn.boxplot()` function requires less boilerplate code, especially when working directly with [pandas DataFrames](#), which is the standard data structure for analytical tasks in Python.

The primary advantage of using a library like [Seaborn](#) for this task is its inherent awareness of data structures, simplifying the mapping of variables to visual elements. Instead of manually preparing lists of values corresponding to each group, as required by the native [Matplotlib](#) function, [Seaborn](#) allows the user to specify the column names for the categorical variable (x-axis) and the quantitative variable (y-axis) directly from the DataFrame. This approach significantly enhances code readability and reduces the potential for errors during data preparation, making it the preferred method for generating complex, grouped visualizations quickly and accurately.

Understanding the Core Syntax for Grouped Boxplots

To generate a grouped [boxplot](#) efficiently, we rely on the integration of [Matplotlib](#) and [Seaborn](#). The fundamental approach involves importing both libraries and then calling the `seaborn.boxplot()`

function, passing in the categorical grouping variable, the numerical data variable, and the [pandas DataFrame](#) containing the required columns. This syntax is highly intuitive and follows a declarative pattern, where the user declares what variables should be plotted rather than how the plot should be constructed step-by-step.

The necessary imports are standard practice in any statistical Python environment. We import [Matplotlib](#) (usually as `plt`) to ensure we have access to the underlying figure and axis controls, and [Seaborn](#) (as `sns`) to utilize its advanced statistical plotting functions. The simplicity of the core command hides its powerful statistical capabilities, automatically calculating the quantiles and identifying potential [outliers](#) for each subgroup defined by the 'group' variable.

The following syntax demonstrates the minimum required parameters to achieve a grouped [boxplot](#) using [Seaborn](#). This assumes your data is stored in a [pandas DataFrame](#) named `df` that contains columns representing the grouping variable (e.g., 'group') and the quantitative measure (e.g., 'values').

You can use the following syntax to create boxplots by group in Matplotlib, leveraging the [Seaborn](#) interface:

```
import matplotlib as plt  
import seaborn as sns
```

```
sns.boxplot(x='group', y='values', data=df)
```

This fundamental command handles all the underlying calculations and rendering, providing a visually appealing and statistically sound [boxplot](#) instantly. The efficiency of this method makes it highly suitable for rapid [data visualization](#) and initial exploratory analysis, allowing the user to focus on interpreting the distributions rather than managing plotting parameters.

Preparing Your Data: Long-Form vs. Wide-Form

When working with statistical plots, the structure of your data is paramount. Data can generally be stored in two primary formats: [long-form data](#) (often called 'tidy' data) and [wide-form data](#). [Seaborn](#), and most modern statistical packages, prefer the [long-form data](#) structure, where each row represents a single observation, and variables are stored in columns. Specifically for a grouped boxplot, this means having one column dedicated to the categorical group identifier (e.g., 'Team') and a second column dedicated to the measured numerical values (e.g., 'Points').

In contrast, [wide-form data](#) structures store different groups or variables across separate columns. For example, if you are measuring points for three teams (A, B, C), the wide format would have three columns, one for 'A', one for 'B', and one for 'C', each containing the point values. While this

format can be easier for manual data entry or basic spreadsheet use, it is suboptimal for plotting functions in Seaborn, which expect defined 'x' and 'y' variables corresponding to individual columns.

Therefore, understanding the format of your input pandas DataFrame is the first critical step. If the data is already in long-form data, the plotting command is straightforward, as shown in Example 1. If the data is in wide-form data, an intermediate step using the `pandas.melt()` function is necessary to transform the structure into the tidy format required by Seaborn's statistical plotting API, as demonstrated in Example 2. This transformation step ensures compatibility and unlocks the full potential of Seaborn for grouped boxplots.

The following examples show how to use this syntax to create boxplots by group for datasets in both long-form and wide-form, addressing the necessary data preparation steps for each structure.

Example 1: Creating Boxplots from Long-Form Data

The following code demonstrates the standard procedure for creating grouped boxplots when the source data is provided in a pandas DataFrame structured in a long-form data (tidy) format. This format is ideal for Seaborn, as the categorical grouping variable ('team') and the numerical measurement ('points') are already distinct columns.

In this example, we utilize NumPy to efficiently create the repeated categorical labels for the teams, ensuring that five data points correspond to each of the three teams (A, B, and C). This setup mimics a dataset where repeated measurements are taken across different groups. The output shows the resulting DataFrame, which is organized perfectly for statistical analysis and visualization, with 15 rows in total, each representing a singular observation of 'points' belonging to a specific 'team'.

The plotting step is remarkably simple: we specify 'team' for the x-axis (the categorical groups) and 'points' for the y-axis (the quantitative values). Seaborn automatically handles the grouping, aggregation, and rendering. We do not need to perform any manual aggregation or reshaping of the data before plotting, making this the most efficient method when data is already tidy. This approach underscores why long-form data is the preferred structure for complex data visualization in Python.

```
import pandas as pd
import numpy as np
import matplotlib as plt
import seaborn as sns
```

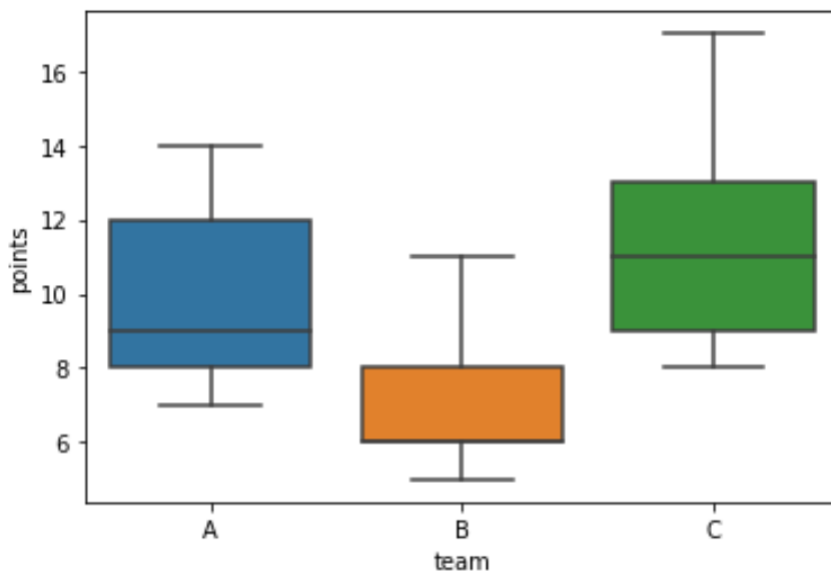
```
#create long-form data
```

```
df = pd.DataFrame({'points': ,
'team': np.repeat(, 5)})

#view data
print(df)

points team
0 7 A
1 8 A
2 9 A
3 12 A
4 14 A
5 5 B
6 6 B
7 6 B
8 8 B
9 11 B
10 8 C
11 9 C
12 11 C
13 13 C
14 17 C

#create boxplot by group
sns.boxplot(x='team', y='points', data=df)
```



Detailed Analysis of Long-Form Plotting Code

The code presented in Example 1 illustrates the power and simplicity of [Seaborn](#) when applied to [long-form data](#). The initial setup involves creating a [pandas DataFrame](#) where the numerical column, 'points', contains all observations, and the categorical column, 'team', defines the grouping structure. The use of the [NumPy](#) function `np.repeat()` is an efficient way to generate the necessary categorical labels for the teams, ensuring data integrity during the creation phase.

The plotting function, `sns.boxplot(x='team', y='points', data=df)`, is the core component. The arguments `x` and `y` directly map the columns in the DataFrame to the axes of the plot. By assigning the categorical variable 'team' to the x-axis, we instruct [Seaborn](#) to calculate a separate [boxplot](#) for all 'points' observations corresponding to each unique value in the 'team' column. This internal mechanism automates the grouping and statistical calculation, which would require manual splitting and aggregation if using base [Matplotlib](#) functions.

Interpreting the output image reveals crucial statistical insights: we can immediately compare the median (the line inside the box) and the interquartile range (the length of the box) for Team A, Team B, and Team C. For instance, if Team C has a higher median and a larger interquartile range than Team B, this visual difference tells us that Team C generally performs higher and has more variability in its scores. This rapid comparison is the essential goal of [grouped boxplots](#), and the tidy data format combined with [Seaborn](#) makes achieving this goal trivial.

Example 2: Creating Boxplots from Wide-Form Data

The following code illustrates the necessary data transformation step required to create grouped

boxplots when the data is structured in a wide-form data format. In this structure, each team's scores are housed in separate columns (A, B, and C). Since Seaborn's `boxplot()` function is designed to work with tidy data where groups are defined by a single categorical column, we must first reshape the data using the powerful pandas.melt() function.

The initial pandas DataFrame setup here demonstrates the typical appearance of wide data: three separate columns holding the numerical measurements. If we attempted to plot this directly using `sns.boxplot()`, the function would not correctly interpret the data as grouped observations. The transformation step is therefore crucial for compliance with the tidy data principles that underpin Seaborn's visualization philosophy. The `melt()` operation transforms the column headers (A, B, C) into values within a new categorical column, while stacking the corresponding numerical values into a single measurement column.

Following the `melt()` operation, the plotting command looks structurally identical to Example 1, demonstrating that once the data is transformed into the required long-form data, the visualization process becomes standardized. However, because the column names generated by `melt()` are generic ('variable' for the group column and 'value' for the measurement column), we must explicitly rename the axis labels using the `.set()` method to ensure the plot is accurately labeled as 'team' and 'points', maintaining clarity for the audience.

```
import pandas as pd
import numpy as np
import matplotlib as plt
import seaborn as sns
```

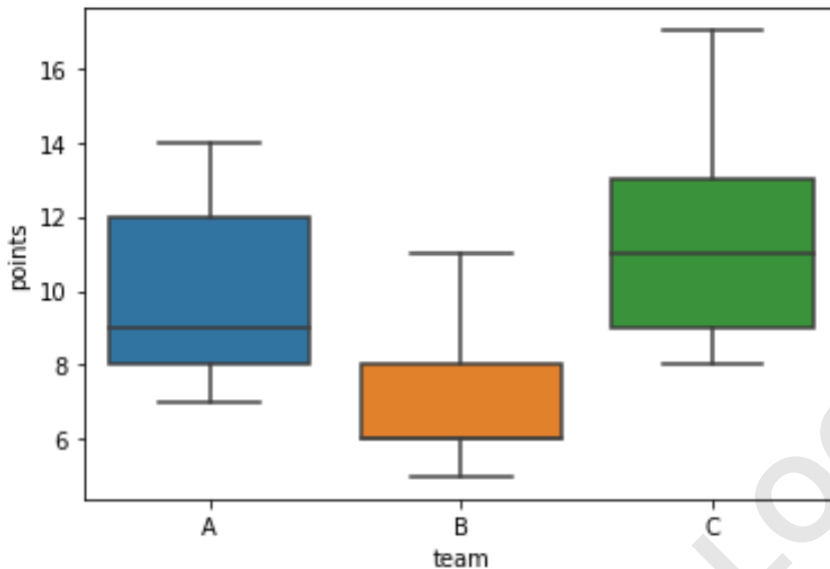
```
#create wide-form data
df = pd.DataFrame({'A': ,
                  'B': ,
                  'C': })
```

```
#view data
print(df)
```

```
A B C
0 7 5 8
1 8 6 9
2 9 6 11
3 12 8 13
4 14 11 17
```

```
#create boxplot by group
```

```
sns.boxplot(x='variable', y='value', data=pd.melt(df)).set(
xlabel='team',
ylabel='points')
```



Handling Wide-Form Data with `pandas.melt()`

The transformation of wide-form data into long-form data is executed by the `pandas.melt()` function, which is essential when data is structured with groups spread across columns. The purpose of melting is to "unpivot" the DataFrame from a wide format to a narrow, or long, format. When called on the DataFrame `df` containing columns A, B, and C, `pandas.melt()` takes these columns and converts them into rows, creating two new columns by default: `'variable'` (holding the original column names, A, B, or C) and `'value'` (holding the corresponding numerical data points).

The resulting melted DataFrame is now in the preferred tidy format, suitable for direct input into Seaborn's statistical plotting functions. The `'variable'` column serves as the categorical grouping variable for the x-axis, and the `'value'` column serves as the quantitative measure for the y-axis. It is important to note that since the `melt()` function is applied directly within the `sns.boxplot()` call--`data=pd.melt(df)`--the transformation is temporary for the purpose of plotting, streamlining the overall code structure.

Following the plotting command, the `.set()` method is chained to the Seaborn call. This is a common practice in Matplotlib/Seaborn to customize the resulting axes properties, such as labels and titles. Since `melt()` uses `'variable'` and `'value'` as default column names, we use `.set(xlabel='team', ylabel='points')` to provide meaningful labels that reflect the true nature of the

data, ensuring the final visual output is clear and professional. This step is crucial for high-quality [data visualization](#).

Customization and Best Practices for Boxplots

While the basic grouped [boxplot](#) is highly informative, [Seaborn](#) and [Matplotlib](#) offer extensive customization options to enhance visual clarity and aesthetic appeal. Customization often includes adding color palettes, adjusting whisker length, showing individual data points, or adding specific annotations. For instance, the `palette` argument in `sns.boxplot()` allows the user to apply different color schemes (e.g., 'viridis', 'Set2'), which can significantly improve the differentiation between groups, especially when presenting the results to a non-technical audience.

One common best practice is to superimpose the raw data points onto the [boxplot](#) using a swarm plot or strip plot, which provides a clearer sense of the sample size and density of the data within each group. This can be achieved by calling `sns.swarmplot()` or `sns.stripplot()` after the `sns.boxplot()` command, using the exact same `x`, `y`, and `data` parameters. By setting the color of the swarm plot points to a distinct color (e.g., black) and adjusting their size, we can reveal potential clustering or overlap that the box summary alone might obscure.

Furthermore, managing plot size and saving the final figure are critical steps in the workflow. While [Seaborn](#) plots utilize [Matplotlib](#) figures, controlling the figure size is often done before the plotting command using `plt.figure(figsize=(width, height))`. Finally, the figure should be saved in a high-resolution format (e.g., PNG or SVG) using `plt.savefig('my_grouped_boxplot.png')` to ensure professional quality for publication or reporting. Adhering to these best practices ensures that the grouped [boxplot](#) is not only statistically sound but also visually optimized for communication.