

# Learning to Visualize Data: Creating Boxplots for Multiple Columns in Seaborn

Authored by  
**Mohammed Iooti**

November 1, 2025

## RECOMMENDED CITATION

Mohammed Iooti (2025). *Learning to Visualize Data: Creating Boxplots for Multiple Columns in Seaborn*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=7545>

[Data visualization](#) serves as a cornerstone of modern data analysis, providing immediate and intuitive access to the underlying structure, distribution, and spread of variables within a dataset. When analysts work with complex tabular data structures, often managed using the robust tools provided by the [Pandas DataFrame](#), the need to perform comparative analysis becomes paramount. Specifically, comparing the statistical distributions of multiple columns simultaneously is a requirement for understanding relative performance or variation across different metrics or groups.

This comprehensive guide is designed to walk expert users through the process of utilizing [Seaborn](#), the powerful Python library built for statistical visualization, to generate a single, consolidated plot. This plot will contain individual [boxplots](#) for several columns, allowing for efficient side-by-side comparison of metrics. While this process is conceptually straightforward, it necessitates a crucial preliminary step: reshaping the data structure to align with the expectations of relational plotting libraries like Seaborn.

## The Challenge of Wide-Format Data in Seaborn

When datasets are initially compiled, they frequently adopt a **wide format**. In this structure, each column represents a distinct group, category, or measured variable. While this format is highly intuitive for data entry and initial inspection, it poses a challenge for many advanced statistical plotting functions. Visualizing all distributions at once offers powerful immediate insights into central tendency, variance, and the presence of outliers across all variables simultaneously.

However, the core functionality of [Seaborn](#), especially methods like `sns.boxplot` which are designed for relational plotting, relies fundamentally on the data being organized in a **long format**. In a long-format structure, there is typically one column dedicated to identifying the categorical variable (the group name, often referred to as 'variable'), and a second column dedicated solely to holding the corresponding measured values (the quantitative data, often referred to as 'value').

Consequently, attempting to directly plot a selection of columns from a wide-format [Pandas DataFrame](#) using Seaborn's relational tools will result in errors or unintended visualizations. To successfully proceed with generating comparative boxplots, we must first master the essential preparatory step of transforming--or "melting"--the data from its wide orientation into the necessary long format. This data manipulation technique ensures compatibility and optimal performance with the visualization library.

## Understanding the Required Seaborn Boxplot Syntax

Before we delve into the data transformation, it is essential to internalize the structural requirements of the [Seaborn](#) `boxplot` function. This function assumes the data has already been prepared into the long format described previously. The function is designed to map two distinct

columns within the provided DataFrame: the column defining the categories (which will form the x-axis groups) and the column defining the quantitative measurements (which will define the y-axis distribution).

The generalized syntax for generating this specific type of visualization strictly enforces the mapping of these variables. The following template illustrates the fundamental command needed to initiate the plotting process once the data is correctly structured:

```
sns.boxplot(x='variable', y='value', data=df)
```

In this standard command, the parameter `x` must reference the column containing the categorical labels (e.g., Team A, Team B), while `y` references the column containing the corresponding numerical measurements (e.g., individual points scored). The `data` parameter, as is customary, points to the newly melted [DataFrame](#) containing the long-format structure. Understanding this syntax prior to manipulation is key to achieving the desired comparative visualization.

## Data Preparation: Creating and Inspecting the Wide DataFrame

To provide a concrete illustration of this technique, we will establish a sample dataset. This example simulates the performance of three distinct basketball teams, labeled A, B, and C, where each column records the points scored by players within that team. Crucially, this initial structure is the typical **wide-format** DataFrame, where each team metric occupies its own column.

Our first step involves importing the necessary [Pandas](#) library and constructing this sample data structure. We define the columns A, B, and C with representative point totals to create a realistic scenario for comparison. This initial setup clearly demonstrates the wide format that we will need to transform:

```
import pandas as pd
```

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

```
#view DataFrame  
df
```

```
A B C  
0 5 8 1  
1 7 8 2
```

```
2 7 9 2
3 9 13 4
4 12 15 5
5 12 17 7
```

Our objective remains to generate three discrete [boxplots](#)--one detailing the distribution for each team (A, B, and C)--and present them side-by-side within a single chart. This side-by-side presentation is the most effective way to visually compare the variance, median performance, and overall spread of points scored across the three groups.

## The Essential Transformation: Melting Data from Wide to Long

The central requirement for plotting multiple column distributions using [Seaborn](#) necessitates moving away from the wide format and adopting the long format. This means that the current column headers (A, B, C) must transition into identifying labels within a new column, and their corresponding numerical values must be stacked into a single measurement column. This process is often referred to as "unpivoting."

This critical transformation is executed with maximum efficiency using the highly specialized [Pandas `melt\(\)` function](#). The primary role of `melt()` is to take a wide-format DataFrame and restructure it into a long format, consolidating the separate measurement columns into two primary columns: `variable` (holding the original column names, which are now the categories) and `value` (holding the numerical data).

By applying the `melt()` function directly to our initial DataFrame `df`, we generate the required structure, which we name `df_melted`. Note that in this standard application, we do not need to specify any `id_vars`, as we want all existing columns to be treated as measured variables:

### #melt data frame into long format

```
df_melted = pd.melt(df)
```

```
#view first 10 rows of melted data frame
```

```
df_melted.head(10)
```

```
variable value
```

```
0 A 5
```

```
1 A 7
```

```
2 A 7
```

```
3 A 9
```

```
4 A 12
```

```
5 A 12
```

6 B 8

7 B 8

8 B 9

9 B 13

The resultant `df_melted` DataFrame is now perfectly structured. It features the column `variable`, which correctly identifies the respective team (A, B, or C), and the column `value`, which contains the associated points scored. This long format ensures that [Seaborn](#) can correctly map the categorical groups against the quantitative values, allowing the visualization to proceed without any further data manipulation.

## Generating the Comparative Boxplot Visualization

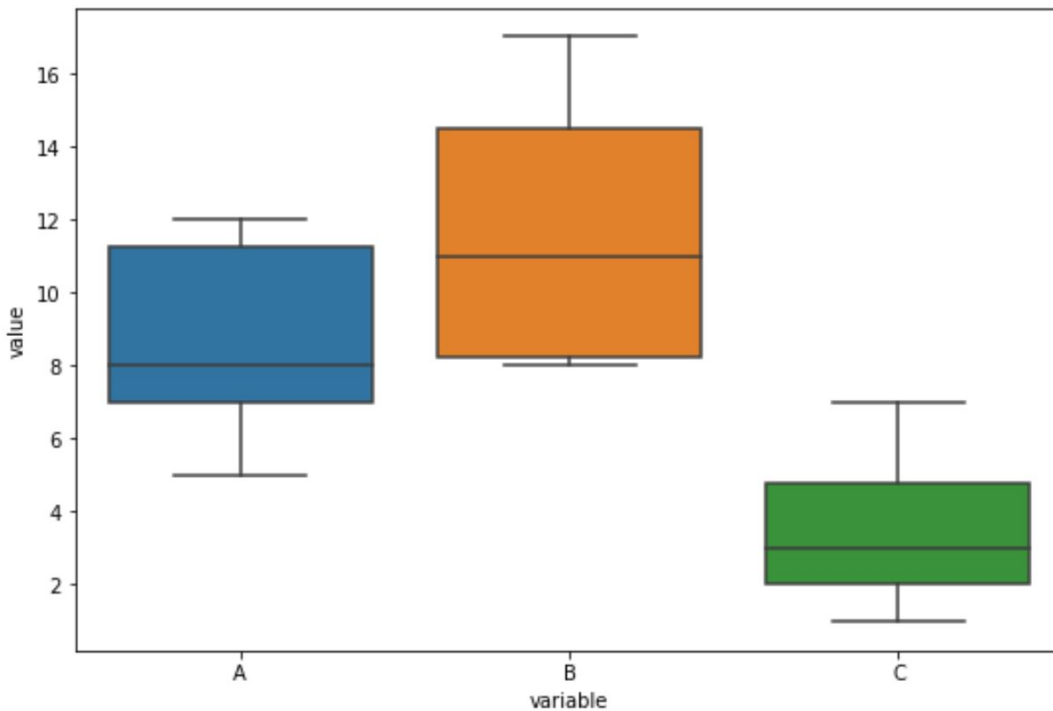
With the data successfully transformed into the long format (`df_melted`), the final step involves generating and customizing the comparative [boxplot](#). This requires ensuring that both [Seaborn](#) and its underlying dependency, [Matplotlib](#) (specifically the `pyplot` module, which is crucial for fine-tuning titles and labels), are correctly imported into the environment.

We execute the plotting command by passing the `df_melted` DataFrame to the `sns.boxplot()` function, explicitly mapping the variables: `x='variable'` handles the team categories, and `y='value'` handles the numerical measurements of points scored. This is where the preparation work pays off:

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

#create seaborn boxplots by group
sns.boxplot(x='variable', y='value', data=df_melted)
```

The initial resulting visualization immediately provides a clear contrast of the points distribution across Teams A, B, and C. The x-axis neatly segregates the teams (the categorical variable), while the y-axis accurately represents the statistical spread and distribution of the quantitative points scored.



## Customizing and Finalizing the Plot for Presentation

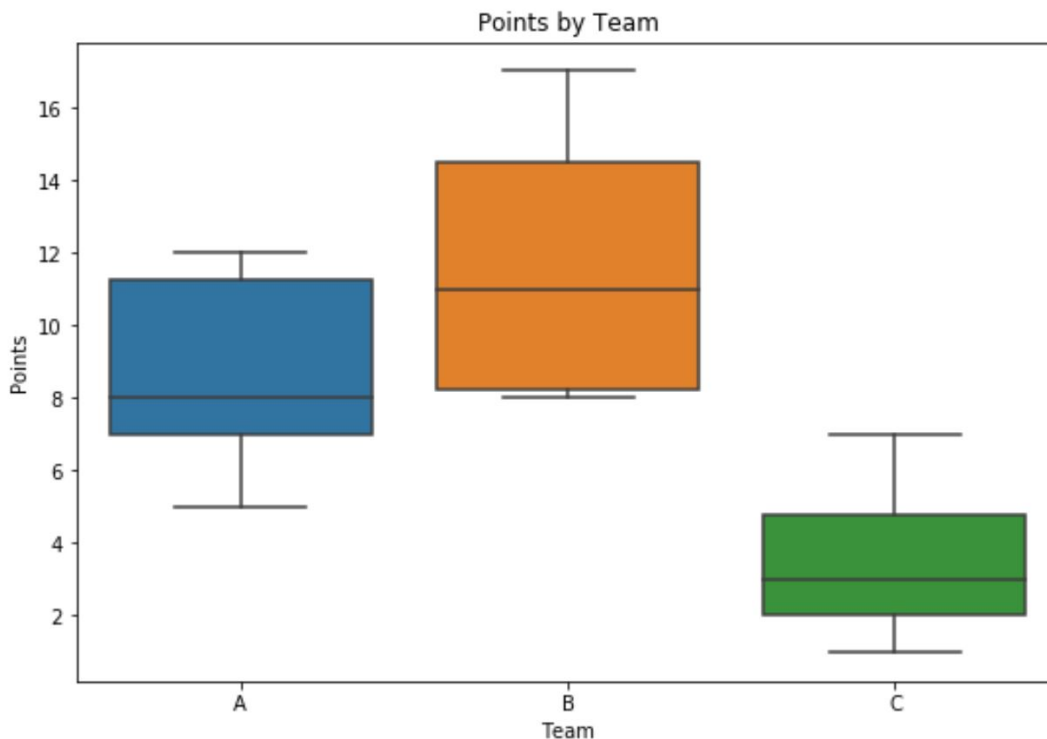
While the initial plot is functionally correct, professional and publication-ready visualizations require clear titles and descriptive axis labels that move beyond the default column names (`variable` and `value`). We leverage methods available within both the [Seaborn](#) plotting object and the imported [Matplotlib](#) module to enhance the chart's comprehensibility and presentation quality.

The following refined syntax demonstrates how to assign a descriptive title ("Points by Team") to the plot using the `set()` method chained onto the `sns.boxplot()` call. Furthermore, we employ the [Matplotlib](#) `plt.xlabel()` and `plt.ylabel()` functions to rename the axes from the generic DataFrame column names to meaningful descriptive labels: 'Team' and 'Points'.

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

```
#create seaborn boxplots by group and set title
sns.boxplot(x='variable', y='value', data=df_melted).set(title='Points by Team')
```

```
#modify axis labels using Matplotlib
plt.xlabel('Team')
plt.ylabel('Points')
```



This customized output represents the final, polished visualization. It provides a statistically robust and visually clear comparison of the performance distributions across the three distinct basketball teams, ready for integration into reports or presentations.

## Summary: Mastering Data Reshaping for Comparative Plots

The process of creating comparative [boxplots](#) for multiple columns using the [Seaborn](#) library is fundamentally dependent on proper data preparation. The core concept mastered in this tutorial is the necessity of transforming a **wide-format** [Pandas DataFrame](#) into a **long-format** structure. This is accomplished efficiently and reliably through the application of the `melt()` function.

By internalizing the requirement of the long format for relational visualizations, users can unlock the advanced capabilities of statistical plotting libraries. This mastery allows for the creation of concise, meaningful visual comparisons of distributions, statistical variance, and central tendencies across numerous variables simultaneously, providing deeper analytical insights.

For users looking to expand their knowledge, we highly recommend consulting the official documentation for further customization options, including adjusting color palettes, managing legend placement, and techniques for efficiently handling much larger datasets that follow similar distributional patterns.

## Additional Resources

[Official Pandas Documentation](#)

[Official Seaborn Documentation](#)

[Matplotlib Pyplot Guide](#)