

Create a Histogram from Pandas DataFrame

Authored by
Mohammed loot

November 3, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Create a Histogram from Pandas DataFrame*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=9079>

Effective [data visualization](#) serves as the cornerstone of exploratory data analysis (EDA), providing analysts with an immediate and intuitive grasp of the underlying distribution of numerical features. Central to this process is the **histogram**, a statistical tool that maps data frequency across defined intervals. This comprehensive guide is designed for Python users, detailing exactly how to leverage the robust capabilities of the [Pandas DataFrame](#) structure to generate informative and highly customizable histograms with remarkable efficiency.

The Foundation: Using Pandas' Built-in `.hist()` Method

The **Pandas** library is engineered for seamless integration with powerful visualization backends, most notably Matplotlib. This integration allows users to generate complex plots using simple, direct methods applied to the data structure itself. The primary tool for distribution analysis is the `.hist()` method, which, when called directly on a **DataFrame** object, offers the quickest and most efficient route to visualizing the shape of a specified numerical column.

The core syntax required to invoke this function is straightforward. It primarily requires specifying the name of the column whose distribution you wish to plot. By minimizing the required arguments for a basic plot, Pandas ensures that preliminary exploratory analysis can be executed with minimal boilerplate code, streamlining the initial stages of any data project.

```
df.hist(column='col_name')
```

The following sections move from theory to application, demonstrating practical examples. We begin with the fundamental task of visualizing a single variable's distribution before exploring techniques that allow for the side-by-side comparison of multiple data distributions.

Implementing the Basic Histogram (Example 1)

To illustrate the fundamental process, we must first establish a sample dataset. For this demonstration, we construct a sample **Pandas DataFrame** containing fictional sports statistics, including points scored, assists, and rebounds. Our objective is to generate an initial [histogram](#) to analyze the spread, shape, and central tendency of a single variable: the `'points'` column.

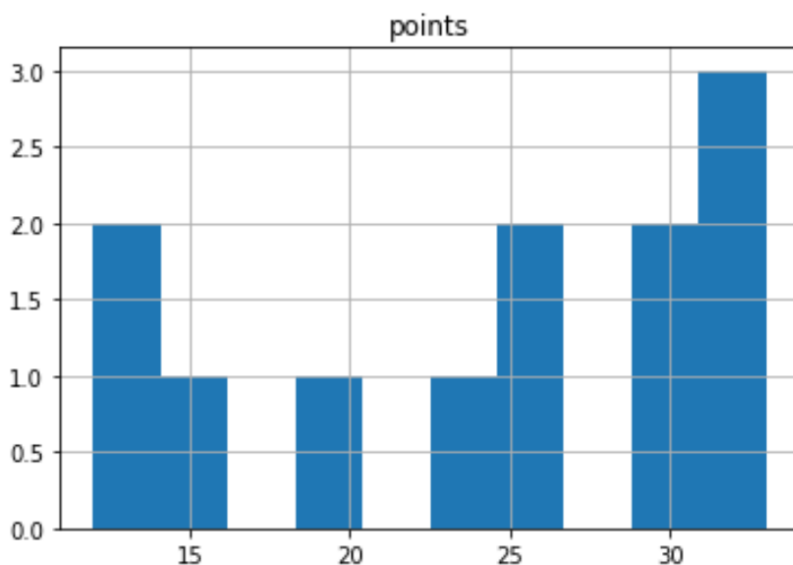
The Python code block below initializes our DataFrame using the standard `pd.DataFrame()` constructor. Subsequently, we call the `.hist()` method, explicitly targeting the `'points'` column. By omitting any additional parameters, we rely on the default visualization settings provided by Pandas and Matplotlib, which are suitable for a quick preliminary view.

```
import pandas as pd
```

```
#create DataFrame
```

```
df = pd.DataFrame({'points': ,  
'assists': ,  
'rebounds': })  
  
#view first five rows of DataFrame  
df.head()  
  
points assists rebounds  
0 25 5 11  
1 12 7 8  
2 15 7 10  
3 14 9 6  
4 19 12 6  
  
#create histogram for 'points' column  
df.hist(column='points')
```

This resulting visualization immediately provides a clear understanding of the data's structure. Analysts can quickly observe where the majority of data points cluster, revealing key characteristics such as symmetry, skewness, or whether the distribution is unimodal or multimodal. This foundational plotting step is crucial before moving on to more complex statistical modeling.



Advanced Customization for Professional Output

While the default histogram provides a good statistical representation, true professional-grade

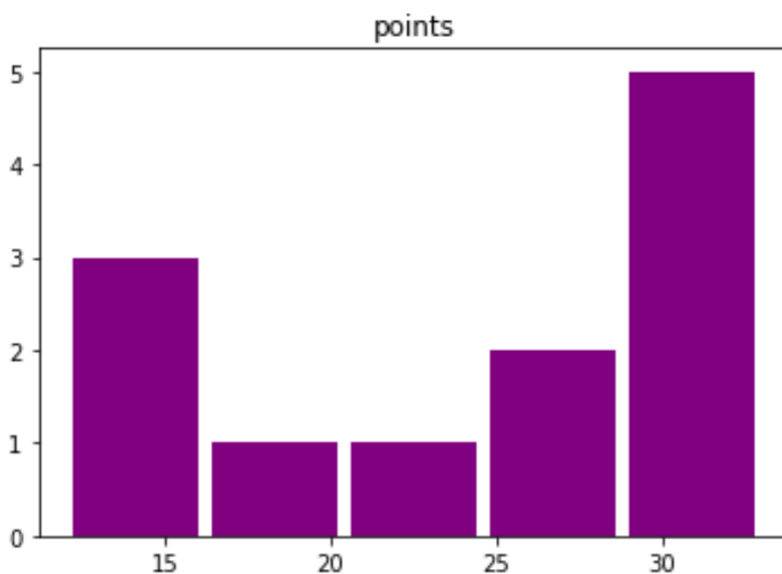
visualizations often require customization to enhance readability and aesthetic appeal. The Pandas `.hist()` method is highly flexible, inheriting numerous arguments from the underlying Matplotlib library, allowing users to fine-tune visual elements such as the number of **bins**, the color scheme, grid visibility, and the width of the bars.

Customization is not merely cosmetic; it is essential for generating high-quality visuals suitable for formal reports or presentations. For instance, adjusting the number of [bins](#) can smooth out or alternatively detail the appearance of the distribution. Disabling the background grid often results in a cleaner, less cluttered look, while specifying a custom color palette ensures compliance with corporate branding or improves visual differentiation.

The following code block demonstrates how to modify several key parameters simultaneously to create a more polished visual representation of the 'points' distribution. We set the number of bins to five, remove the background grid, slightly narrow the bars using `rwidth`, and apply a specific color.

```
#create custom histogram for 'points' column  
df.hist(column='points', bins=5, grid=False, rwidth=.9, color='purple')
```

In this resulting custom plot, the **x-axis** effectively displays the range of points scored per player, grouped into the five specified [bins](#) or intervals. Crucially, the **y-axis** now indicates the [frequency](#), representing the total count of players whose scores fall within each defined score bracket. This presentation clearly communicates the density of observations across the measured range.



Comparative Analysis: Grouping Data with the by Parameter (Example 2)

One of the most powerful features of the Pandas `.hist()` method is its capacity to facilitate comparative analysis by generating multiple visualizations based on a categorical grouping variable. This functionality is enabled by the `by` argument, which automatically splits the underlying numerical data and plots a separate **histogram** for every unique value present in the specified grouping column. This feature is an indispensable tool for comparing distributions across distinct segments or groups within a larger dataset.

To demonstrate this technique, we construct a new [DataFrame](#) that explicitly incorporates a categorical `'team'` variable alongside the numerical `'points'` data. Our objective is to instruct Pandas to plot the distribution of `'points'` separately for Team A and Team B. This setup allows for an immediate, side-by-side comparison of the scoring patterns and typical performance levels of the two teams.

The following code defines the new dataset and then executes the `.hist()` call, utilizing the `by='team'` parameter. Notice that we retain the customization arguments (such as `bins`, `grid`, and `rwidth`) for visual clarity, while introducing `sharex=True`, a crucial element for accurate comparison, which will be discussed next.

import pandas as pd

```
#create DataFrame
```

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

```
#view first five rows
```

```
df.head()
```

```
team points
```

```
0 A 25
```

```
1 A 12
```

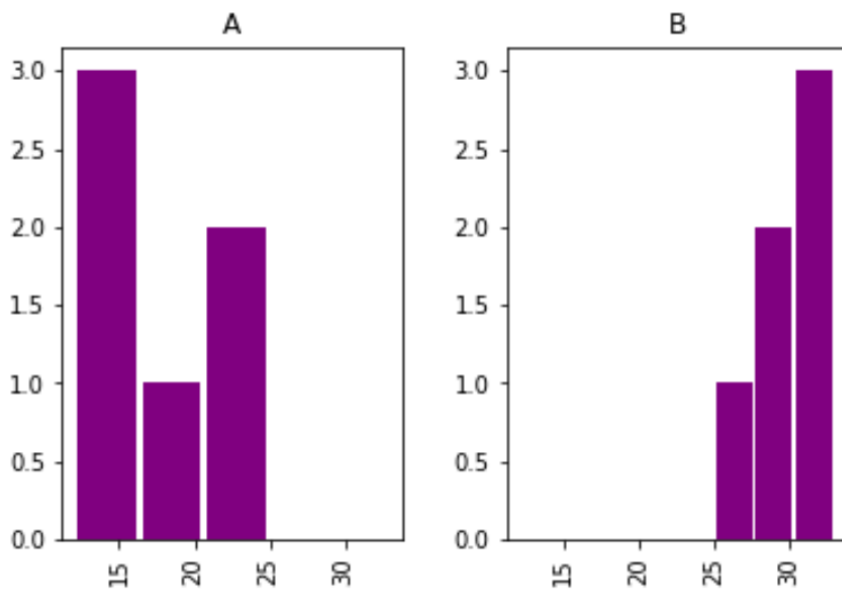
```
2 A 15
```

```
3 A 14
```

```
4 A 19
```

```
#create histogram for each team
```

```
df.hist(column='points', by='team', bins=3, grid=False, rwidth=.9,  
color='purple', sharex=True)
```



Ensuring Accuracy: Controlling Axis Scaling

When generating multiple subplots for comparative analysis, as done using the `by` parameter, maintaining visual consistency is paramount for drawing accurate conclusions. The `sharex` argument within the `.hist()` method directly addresses the potential for visual distortion that can occur if subplots automatically adjust their ranges based solely on local data.

By setting the `sharex=True` parameter in the previous code example, we explicitly mandate that both generated [histograms](#) must utilize the exact same scale and limits for the x-axis (the variable being measured, 'points'). If this parameter were omitted, one team's plot might appear compressed or expanded relative to the other, potentially misleading the viewer regarding the true differences in their score distributions.

This shared axis specification significantly enhances the viewer's ability to compare the distribution of values between the two groups objectively. Analysts frequently rely on this comparative approach to identify significant disparities in data performance, skewness, or central tendencies across defined categories, making `sharex` a mandatory setting for effective comparative data visualization.

Mastering Key Parameters of the `.hist()` Function

To fully harness the power of histogram creation in Pandas, it is essential to possess a solid understanding of the primary parameters available within the `.hist()` function. These parameters collectively give the user granular control over both the statistical grouping of the data and the final aesthetic output of the visualization.

column: This mandatory parameter defines the numerical column(s) from the **DataFrame** that the user intends to plot. It accepts either a single string (for one histogram) or a list of strings, which instructs Pandas to plot a separate histogram for each variable listed.

by: Specifies a categorical column used to group the data. This parameter is integral to comparative analysis (as shown in Example 2), as it generates separate subplots for each unique category found within the grouping column.

bins: Determines the number of [bins](#) (intervals) used to divide the entire data range. The choice of bin count is a critical decision in statistical visualization: a high number provides greater detail but can introduce noise, while a low number smooths the distribution but risks hiding subtle features. Selecting the optimal number of bins is often an iterative process.

grid: A boolean value (True or False) that dictates whether the background grid lines are displayed. Setting it to `False`, as demonstrated in our examples, is a common practice for achieving a cleaner, more professional visual appearance.

rwidth: Controls the relative width of the histogram bars compared to the width of the bin interval. Setting `rwidth` slightly below 1.0 (e.g., 0.9) introduces small gaps between the bars, significantly improving visual separation and making the delineation between columns clearer.

sharex/sharey: These boolean parameters control whether subplots share the same x or y axes, respectively. As established, `sharex=True` is essential for ensuring accurate, non-misleading comparison in plots generated using the `by` parameter.

color: Allows the user to specify the fill color of the histogram bars using a recognized string (e.g., 'purple', 'blue') or a precise hex code, enabling aesthetic control.

Conclusion and Further Visualization Resources

Generating a [histogram](#) directly from a [Pandas DataFrame](#) via the concise `.hist()` method is a foundational and indispensable technique for understanding data distribution in Python. By thoughtfully utilizing key parameters such as `column`, `by`, and `bins`, data analysts can efficiently transition from raw numerical data to insightful visualizations. This mastery allows for quick examination of a single variable's [frequency](#) profile or the simultaneous comparison of distributions across multiple defined groups. Proficiency in these basic plotting functions forms the essential skill set required for more advanced [data visualization](#) and statistical analysis within the Python ecosystem.

For those interested in expanding their visualization toolkit beyond **histograms**, the following tutorials explain how to create other common plots in Python: