

Learn Univariate Analysis with Python: A Beginner's Guide

Authored by
Mohammed loot

November 1, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learn Univariate Analysis with Python: A Beginner's Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=7833>

The concept of [Univariate Analysis](#) is foundational in data science, representing the rigorous examination of a **single variable** within a larger dataset. Derived from the prefix "uni" meaning "one," this methodology exclusively focuses on characterizing one attribute at a time--specifically its distribution, measures of central tendency, and overall dispersion. Univariate analysis is the essential first step in [Exploratory Data Analysis \(EDA\)](#), as it furnishes critical preliminary insights necessary before progressing to more sophisticated statistical modeling or machine learning tasks.

To comprehensively understand a variable's inherent characteristics, data scientists rely on a combination of numerical and graphical methods. These three core approaches provide a holistic view of the data's structure and behavior:

Summary Statistics: This involves calculating quantitative measures such as the mean, median, mode, and standard deviation. These metrics are crucial for rapidly quantifying the center, variability (dispersion), and potential skewness of the variable's values.

Frequency Distributions: Creating tables that count the occurrence of unique values or predefined ranges (bins). This tabular method is invaluable for determining the variable's modality and identifying the most common observations.

Data Visualizations: Employing powerful graphical tools, including histograms, [box plots](#), and density plots, to visually interpret the overall shape, symmetry, and presence of outliers within the distribution.

This comprehensive guide will walk through the execution of these fundamental univariate techniques using the ubiquitous programming language, [Python](#), specifically leveraging its premier data manipulation library, **Pandas**. To provide a practical demonstration, we will analyze a sample [Pandas DataFrame](#) structured to track hypothetical athlete performance metrics across several categories (points, assists, rebounds).

```
import pandas as pd
```

```
#create DataFrame
```

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

```
#view first five rows of DataFrame
```

```
df.head()
```

```
points assists rebounds
```

```
0 1.0 5 11
```

```
1 1.0 7 8
```

```
2 2.0 7 10
```

3 3.5 9 6

4 4.0 12 6

Calculating Descriptive Summary Statistics

The essential entry point for univariate exploration is the computation of descriptive [summary statistics](#). These numerical measures serve as concise descriptors, quantifying fundamental characteristics such as the central location, the degree of dispersion, and the shape of the data distribution. When dealing with continuous variables, such as the 'points' column in our athlete dataset, these statistics allow for the immediate identification of representative or 'typical' values and highlight the inherent variability within the observed scores.

Our focus here will be on three pivotal measures for the 'points' variable: the arithmetic **mean** (the average value), the **median** (the middle value when data is ordered), and the **standard deviation** (which quantifies the typical distance of observations from the mean, acting as a measure of spread). One of the greatest advantages of using the **Pandas** library is the ease with which these complex calculations can be performed; simple, built-in methods are applied directly to the selected DataFrame column.

The following [Python](#) syntax illustrates how to retrieve these individual metrics efficiently using the dedicated Pandas functions:

```
#calculate mean of 'points'
```

```
df.mean()
```

```
5.706666666666667
```

```
#calculate median of 'points'
```

```
df.median()
```

```
5.0
```

```
#calculate standard deviation of 'points'
```

```
df.std()
```

```
3.858287308169384
```

Generating a Frequency Distribution Table

In contrast to aggregate summary statistics that condense data into single metrics, a [frequency table](#) offers a granular, value-by-value perspective on the dataset's composition. This technique is

invaluable, particularly when examining discrete or categorical data, as it explicitly enumerates the rate of occurrence for every distinct measurement. Understanding the frequency of observations helps reveal underlying patterns that might be obscured by the mean or median alone.

When applied to our continuous 'points' variable, the frequency table effectively counts the number of times each unique score appears within the column. This tabulation immediately highlights the distribution's **mode**--the value that registers the highest frequency. In **Pandas**, this functionality is accessed using the highly useful `value_counts()` method, which automatically sorts the results by frequency in descending order.

#create frequency table for 'points'

```
df.value_counts()
```

```
4.0 3
```

```
1.0 2
```

```
5.0 2
```

```
2.0 1
```

```
3.5 1
```

```
6.5 1
```

```
7.0 1
```

```
7.4 1
```

```
8.0 1
```

```
13.0 1
```

```
14.2 1
```

```
Name: points, dtype: int64
```

A quick analysis of the generated frequency table yields several immediate and meaningful insights regarding the athlete performance data. By examining the counts associated with each unique score, we can draw the following conclusions about the distribution:

The score of **4.0** appears 3 times, confirming it as the primary mode of this specific dataset.

The scores **1.0** and **5.0** are tied for the second-most frequent observation, each appearing 2 times.

A significant portion of the data consists of unique values (2.0, 3.5, 6.5, 13.0, etc.), indicating a high level of distinct measurements across the higher end of the scoring range.

Visualizing Distribution with Charts

Although numerical summaries provide precision, the overall shape, symmetry, and presence of extreme outliers in a distribution are best communicated through visual inspection. Consequently, visualization is an indispensable component of effective [univariate analysis](#). Charts allow the data

analyst to immediately detect structural characteristics such as skewness, modality, or unexpected gaps that are easily overlooked when reviewing purely statistical tables.

To demonstrate the power of graphical techniques, we will generate three distinct plot types. These visualizations rely on the robust charting ecosystem of [Python](#), specifically utilizing the foundational plotting library, [Matplotlib](#), and the high-level statistical visualization library, [Seaborn](#).

Method 3a: The Box Plot for Five-Number Summary

The [box plot](#), often referred to as a box-and-whisker plot, is a streamlined visualization tool designed to quickly summarize the distribution based on five key numbers: the minimum value, the first quartile (Q1), the median (Q2), the third quartile (Q3), and the maximum value. Its structure makes it exceptionally effective for rapidly discerning the data's central tendency, spread, and, critically, for flagging potential outliers that fall outside the typical range defined by the whiskers.

Although Box plots are statistical charts, we can generate them directly using the Pandas `boxplot()` method. This function internally relies on [Matplotlib](#) for rendering but simplifies the syntax significantly, allowing us to easily visualize the distribution characteristics for our 'points' variable.

```
import matplotlib.pyplot as plt
```

```
df.boxplot(column=, grid=False, color='black')
```



The resulting box plot provides immediate visual evidence of the data's shape. The central box encapsulates the Interquartile Range (IQR), representing the spread of the middle 50% of the scores. Crucially, the chart suggests that the distribution of 'points' is positively skewed (right-skewed), a conclusion supported by the elongated upper whisker and the presence of data points marked individually above the upper boundary, which represent potential extreme outliers in athlete performance.

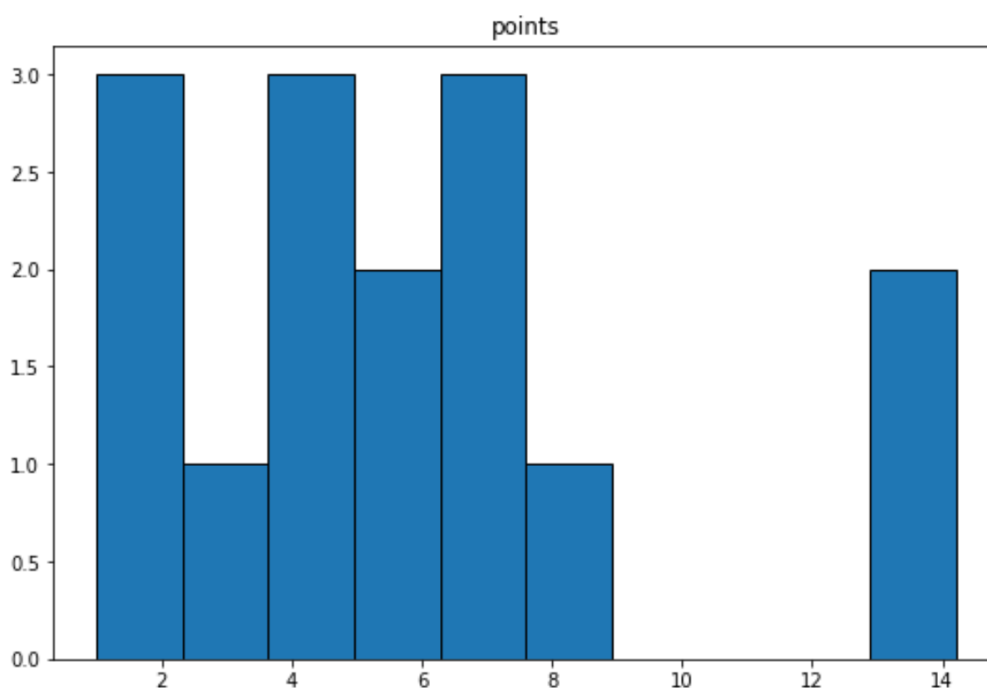
Method 3b: Histogram for Frequency Bins

The [histogram](#) remains the paramount visualization technique in [univariate analysis](#). This chart functions by partitioning the entire range of observed values into contiguous intervals, known as bins. It then plots the frequency (count) of observations that fall within each bin, thereby delivering a clear visual representation of the underlying probability distribution of the variable.

Similar to the box plot, generating a histogram in Python is streamlined using the `hist()` method, which is applied directly to the [Pandas DataFrame](#) object. This convenience allows for rapid statistical plotting without extensive configuration, producing a histogram for the designated variable, 'points', instantly.

```
import matplotlib.pyplot as plt
```

```
df.hist(column='points', grid=False, edgecolor='black')
```



The visual pattern revealed by the histogram strongly corroborates the findings from the box plot. We observe a pronounced concentration of athlete performance data within the lower score bins (specifically 0 to 5 points). Furthermore, the bars diminish gradually as they extend toward the higher point totals, creating a distinct "long tail" to the right. This shape provides definitive confirmation that the variable 'points' exhibits a positively (right) skewed distribution.

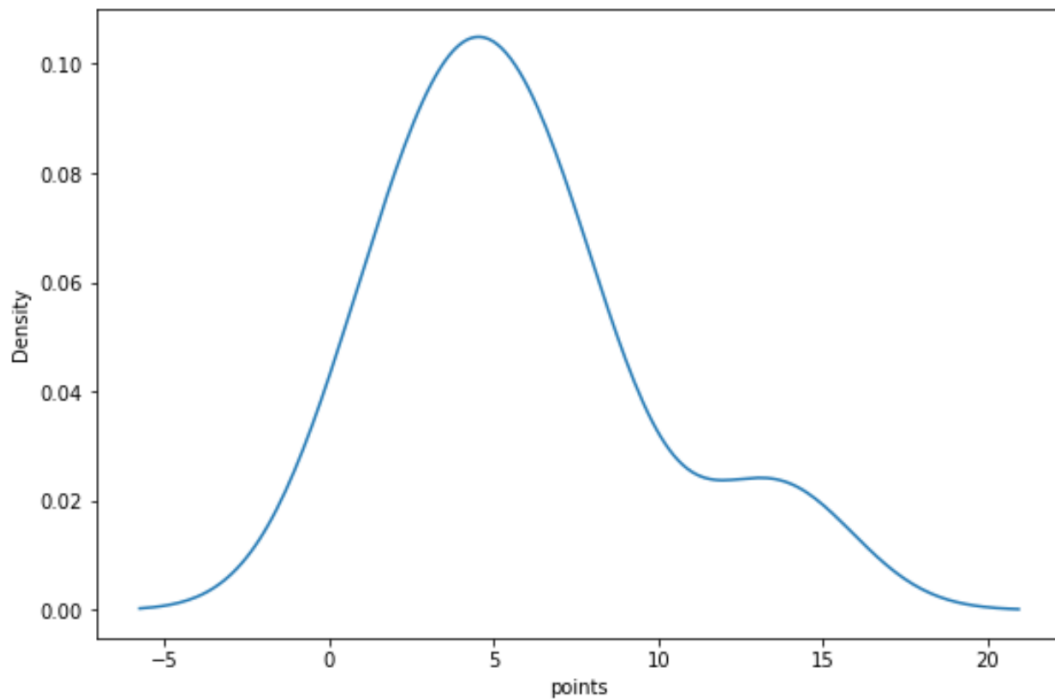
Method 3c: Kernel Density Estimation (KDE) Plot

To achieve a more refined and smoother graphical representation of the distribution, data scientists often employ the [Kernel Density Estimation \(KDE\)](#) plot. The key distinction between a KDE plot and a histogram lies in its methodology: rather than relying on discrete, fixed bins, the KDE plot estimates the underlying probability density function of the random variable, producing a smooth, continuous curve. This results in a cleaner visual analysis of the distribution's overall shape, mitigating the visual artifacts inherent in the binning process.

For generating high-quality statistical plots like the KDE, the **Seaborn** library is highly recommended. Built upon [Matplotlib](#), Seaborn offers streamlined functions for complex visualizations. We use the `kdeplot()` function to render the estimated density curve for the 'points' variable, further elucidating its distribution shape.

```
import seaborn as sns
```

```
sns.kdeplot(df)
```



The resulting KDE plot serves as the clearest visual synthesis of our findings. The curve exhibits a prominent peak, indicating a high probability density centered around the lower scores (1 to 5 points), followed by a gentle, extended tail toward the right. This continuous curve definitively confirms the heavy right-skewness of the distribution. By employing descriptive statistics, frequency tables, and these complementary visualizations, we have successfully performed a comprehensive [Exploratory Data Analysis \(EDA\)](#) on this single variable, 'points', fulfilling the requirements of effective univariate analysis.