

Learning How to Calculate the Median Using Pandas

Authored by
Mohammed loot

November 4, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning How to Calculate the Median Using Pandas*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=9583>

Calculating the [median](#) is a cornerstone technique in [exploratory data analysis](#) (EDA). Serving as a crucial measure of [central tendency](#), the median is defined as the middle value of a dataset when the observations are ordered sequentially. Its primary advantage over the mean (average) lies in its inherent resistance to the distorting effects of statistical **outliers** and extreme values. This robustness makes it the preferred metric in fields ranging from finance to environmental science when summarizing skewed distributions.

When working within the powerful Python ecosystem, especially with structured tabular data managed by the [Pandas](#) library, determining the median is simplified immensely. Pandas provides highly optimized, vectorized functions that allow data professionals to calculate this measure efficiently across vast datasets. This comprehensive guide, tailored for data analysts and Python developers, illustrates the precise methods for applying the built-in `median()` function across various scopes: individual columns, specified subsets of columns, and all quantitative features within a [DataFrame](#) object.

The Statistical Importance of the Median in Data Analysis

Understanding the context behind the median calculation is vital for accurate data interpretation. Unlike the arithmetic mean, which shifts dramatically when influenced by high or low extreme values, the median remains stable because it focuses only on the positional center of the data. For example, when analyzing income distribution--a typically skewed dataset--the median provides a far more representative figure of the 'typical' income than the mean, which can be inflated by a small number of extremely wealthy individuals. Utilizing the median, therefore, ensures a more truthful and less biased summary of the data's central point.

The calculation process involves two main steps: first, sorting the numerical observations from smallest to largest, and second, identifying the middle point. If the number of observations (N) is odd, the median is the single value at the position $(N+1)/2$. If N is even, the median is calculated as the average of the two middle values ($N/2$ and $N/2 + 1$). Pandas handles this underlying logic automatically, allowing analysts to focus purely on the data manipulation and interpretation stages.

Mastering the Pandas `median()` Function

The fundamental mechanism for calculating the median in this environment is the `median()` function, which is tightly integrated into the [Pandas](#) library structure. This function is an aggregation method that operates seamlessly on both [Series](#) (single columns) and [DataFrame](#) objects. Its design prioritizes flexibility and speed, making complex statistical summaries achievable with minimal lines of code.

When applied to a [DataFrame](#), the `median()` function calculates the median column-wise by

default (using `axis=0`). It is also engineered to intelligently handle mixed data types: it automatically excludes non-numeric columns (like strings or timestamps) unless explicitly instructed otherwise, focusing only on quantitative features where a statistical median is meaningful. Furthermore, its inherent handling of missing data, controlled by the `skipna` parameter, is critical for real-world datasets.

The following list details the three most common and practical scenarios for applying the `median()` function, guiding users from targeted calculations to broad data profiling:

Targeted Calculation: Determining the central value within a single, specific column (a Pandas [Series](#)).

Comparative Analysis: Calculating the median across a user-defined subset of numeric columns simultaneously.

Holistic Profiling: Quickly obtaining the median for every quantifiable (numeric) column within the entire DataFrame.

The syntax below provides a quick reference for these essential use cases, demonstrating the versatility of the aggregation method:

Finding the median value in a specific column (Series object)

```
df.median()
```

```
# Finding median values across multiple specified columns
```

```
df[columns].median()
```

```
# Finding the median value across every numeric column in the DataFrame
```

```
df.median()
```

Preparing the Dataset for Practical Demonstration

To fully grasp the mechanics of the `median()` function, we must work with a realistic sample dataset. We will construct a small [DataFrame](#) designed to mimic player statistics. This dataset is intentionally structured to contain a mix of different data types--specifically, a string identifier ('player') and multiple numeric performance metrics ('points', 'assists', 'rebounds'). Crucially, we have also introduced an explicit missing value, represented by `pd.NA`, within the 'points' column. This inclusion allows us to demonstrate how Pandas robustly handles incomplete data during statistical aggregation.

The creation script below sets the foundation for all subsequent examples. By observing the structure, we can clearly anticipate how the `median()` function will differentiate between the categorical data and the quantitative data, ensuring that only the relevant numerical arrays are

processed. This setup is paramount for illustrating the function's default behavior, particularly its ability to ignore non-computable entries.

Import Pandas and create the sample DataFrame

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

# Display the resulting DataFrame structure
df
```

```
player points assists rebounds
0 A 25 5 11
1 B NA 7 8
2 C 15 7 10
3 D 14 9 6
4 E 19 12 6
5 F 23 9 5
6 G 25 9 9
7 H 29 4 12
```

Example 1: Calculating the Median for a Single Pandas Series

The most straightforward application involves extracting the [median](#) for a single variable of interest. This step is essential when performing deep statistical dives into individual metrics. To achieve this, we first use standard bracket notation (`df`) to select the column, which results in a Pandas [Series](#) object. The `median()` method is then applied directly to this Series.

We focus here on the **'points'** column. As highlighted during the setup, this column contains one missing value (`NA`). A key feature of the Pandas `median()` function is its default handling of such data: it automatically invokes the parameter `skipna=True`. This means the calculation will proceed only on the valid, non-missing numeric entries, preventing the missing data from skewing or invalidating the result. The calculation effectively uses $N=7$ observations instead of the total $N=8$ rows.

Executing the calculation for the 'points' column yields the following result:

```
# Calculate the median value of the 'points' column
df.median()
```

23.0

The calculated [median](#) value is **23.0**. To verify this, one would sort the seven valid 'points' observations (14, 15, 19, 23, 25, 25, 29). Since $N=7$ (an odd number), the median is the fourth value in the sorted list, which is 23. This result provides a reliable measure of the central performance level for the players recorded in this dataset.

Example 2: Vectorized Median Calculation Across Specified Columns

In comparative statistics or feature engineering, analysts often need to summarize multiple related variables simultaneously. [Pandas](#) facilitates this high-efficiency calculation through its vectorized operations. By passing a Python list of column names (using double brackets, e.g., `df[]`) to the [DataFrame](#), we select a subset of columns, and applying `.median()` to this subset executes the calculation for each column independently.

In this example, we calculate the median for both the **'points'** and **'rebounds'** metrics. This method returns a new Pandas [Series](#) where the index consists of the column names, providing an instantaneous, neatly summarized table of results. This approach saves significant computational time compared to iterating through columns individually, which is crucial when dealing with datasets containing hundreds of metrics.

The execution and output for the selected columns are as follows:

```
# Calculate median values for 'points' and 'rebounds' columns
df[].median()
```

```
points 23.0
rebounds 8.5
dtype: float64
```

While the median for 'points' remains 23.0, the median for 'rebounds' is calculated as **8.5**. This fractional result is derived from the 'rebounds' column having an even number of observations ($N=8$, with no missing values). Sorting the 'rebounds' data (5, 6, 6, 8, 9, 9, 10, 11, 12) reveals that the two middle values are 8 and 9. The median is thus computed as the arithmetic average of these two central observations: $(8 + 9) / 2 = 8.5$. This demonstrates the dual nature of median calculation based on the parity of the sample size.

Example 3: Comprehensive Median Calculation Across All Numeric Features

During the initial stages of data profiling, analysts often require a rapid summary of all quantitative

features within a dataset. The [median\(\) function](#) is designed to facilitate this by allowing direct application to the entire [DataFrame](#) object without specifying any column names. When executed this way, Pandas intelligently scans the data types of all columns.

The function automatically applies the calculation only to columns designated as numeric (typically integer or float data types) while silently bypassing non-numeric columns, such as the 'player' (string) column in our example. The result is a concise, high-level summary of the central tendency for every quantifiable metric, offering a robust baseline for further statistical modeling or visualization efforts.

Applying the `median()` function globally:

```
# Calculate the median value of all numeric columns in the DataFrame
```

```
df.median()
```

```
points 23.0  
assists 8.0  
rebounds 8.5  
dtype: float64
```

The output provides the median for 'points' (23.0), 'assists' (8.0), and 'rebounds' (8.5), offering a comprehensive snapshot of the central tendency for all quantitative measures in our sample dataset. The 'assists' column has eight observations, and when sorted (4, 5, 7, 7, 9, 9, 9, 12), the two middle values are 7 and 9. The median is thus calculated as $(7 + 9) / 2 = 8.0$.

Advanced Handling of Missing Data (`skipna` Parameter)

One of the most powerful and often overlooked aspects of the Pandas aggregation functions, including [median\(\)](#), is its default handling of missing values. As observed in the 'points' column, the function successfully calculated the [median](#) (23.0) despite the presence of a missing entry (`pd.NA`). This is achieved because the `skipna` parameter is set to `True` by default.

Setting `skipna=True` instructs Pandas to exclude any entries identified as missing (e.g., `NaN`, `None`, or `pd.NA`) from the calculation before attempting to find the middle value. This ensures that the calculated median is a true representation of the valid observed data points, enhancing the reliability of the [central tendency](#) measure. This default setting is crucial for maintaining data integrity in real-world scenarios where data collection is rarely perfect.

However, analysts occasionally need to enforce strict handling of missing data. If the user explicitly sets `skipna=False`, the presence of even a single missing value within the array will cause the aggregation function to return `NaN` (Not a Number) for that calculation. This setting is useful as a

diagnostic tool to flag columns that contain gaps, forcing the analyst to address data imputation or cleaning before proceeding with statistical summaries. For instance, recalculating the median of 'points' with `skipna=False` would return `NaN` because the column contains `pd.NA`.

Further Exploration and Related Statistical Methods

Mastering the calculation of the median is just one step in becoming proficient in statistical data analysis using Python and [Pandas](#). To deepen your understanding of quantitative analysis and data manipulation, consider exploring related aggregation methods and key parameters that govern how these statistics are computed.

The following topics are highly recommended for continued learning and represent essential skills for any data professional:

Calculating Other Measures of Central Tendency: Explore the complementary methods for determining the [mean](#) (`.mean()`) and the [mode](#) (`.mode()`). Understanding how these three measures differ, particularly in skewed datasets, is fundamental to robust statistical reporting.

Understanding the `axis` Parameter: Learn how to control the direction of aggregation. While `axis=0` (column-wise) is the default and most common setting, specifying `axis=1` enables row-wise calculations, allowing you to compute the median across metrics for each individual observation.

Advanced Missing Data Strategies: Beyond simply skipping missing values, explore techniques for data imputation (filling missing values with estimates, such as the calculated median itself) using methods like `.fillna()`.

Dispersion Metrics: Complement your central tendency measures by calculating metrics of data spread, such as standard deviation (`.std()`) and variance (`.var()`), to gain a complete picture of the dataset's distribution.

By effectively utilizing the robust tools provided by the Pandas library, data analysts can quickly and accurately derive meaningful statistical insights, ensuring that foundational metrics like the median are calculated efficiently and reliably.