

# How to Calculate the Median Using Pandas: A Step-by-Step Guide

Authored by  
**stats writer**

December 4, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Calculate the Median Using Pandas: A Step-by-Step Guide*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=105293>

The Pandas library is arguably the most essential tool for Python-based data science and analysis. It provides robust data structures, most notably the DataFrame, and comprehensive functions for data manipulation and statistical evaluation. When working with numerical datasets, understanding the central tendency is critical, and the median offers a highly reliable measure, particularly in the presence of outliers. This comprehensive guide details the efficient methods provided by Pandas to calculate the median value of your data.

To calculate the median of a set of values in Pandas, you leverage the powerful built-in method, median(). This function operates on a Series (a single column) or an entire DataFrame, returning the middle value of the sorted data distribution. Unlike the mean, the median is resistant to extreme values, making it an indispensable statistic for summarizing skewed distributions. We will explore how to apply this function across single columns, multiple selected columns, and all available numerical features within a dataset.

Furthermore, while the primary focus is on the **median()** function, it is important to remember that Pandas offers broader statistical capabilities. Tools like the **describe()** function provide a quick overview of various descriptive statistics, including the mean, standard deviation, minimum, maximum, and quartiles, alongside the median. By mastering the application of the **median()** method, data analysts can quickly and accurately derive meaningful insights from their datasets, paving the way for more rigorous statistical modeling.

## The Core Function: Using `median()` in Pandas

The Pandas library simplifies statistical calculation through vectorized operations. The fundamental way to find the middle value in a dataset is by invoking the **median()** function directly on a Pandas Series or DataFrame. This method is highly flexible, allowing analysts to target specific subsets of the data based on their analytical needs. Whether you need the central value for a single feature or a comprehensive summary across many variables, the syntax remains intuitive and powerful.

Understanding the application scope is key. When applied to a single column (Series), **median()** returns a scalar value. When applied to a DataFrame containing multiple numeric columns, it returns a Series where the index labels are the column names and the values are the corresponding medians. By default, it operates column-wise, meaning it calculates the median for each column independently. This efficiency makes quick statistical summaries trivial, even for very large datasets.

The following structures demonstrate the primary ways to utilize the **median()** function to find the median of one or more columns within a Pandas DataFrame, illustrating its versatility:

### # Find the median value in a specific, single column

## **df.median()**

```
# Find the median value across several specified columns, returning a Series of results
```

```
df].median()
```

```
# Find the median value in every numeric column within the entire DataFrame
```

```
df.median()
```

## **Setting Up the Example Dataset (Pandas DataFrame)**

To demonstrate the practical application of the **median()** function, we will establish a sample Pandas DataFrame representing player statistics. This dataset includes a mix of categorical data (player name) and several numeric columns (points, assists, rebounds). Crucially, this dataset is designed to include a missing value (**pd.NA**) in the 'points' column, allowing us to observe how Pandas handles incomplete data during median calculation.

The creation of a well-defined example allows us to isolate the behavior of the statistical functions. The DataFrame is instantiated using a dictionary structure, where keys correspond to column headers and values are lists of data points. Observing the initial structure, including the specific data types, is important before proceeding with any statistical aggregation, as the **median()** function will only operate on appropriate numeric types, silently skipping non-numeric columns like 'player'.

The following code block executes the creation of the sample DataFrame and then displays its contents, which serves as the foundation for all subsequent examples in this guide:

```
# Create the sample DataFrame using Python and Pandas functions
```

```
df = pd.DataFrame({'player': ,
```

```
'points': ,
```

```
'assists': ,
```

```
'rebounds': })
```

```
# View 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: Finding the Median of a Single Column

Calculating the median for a single, focused variable is the most straightforward application of the Pandas method. This is achieved by first selecting the desired column--which returns a Pandas Series--and then chaining the **median()** function directly onto that Series. For instance, if an analyst is specifically interested in the central performance metric for 'points', this targeted approach provides the precise statistical measure without calculating unnecessary statistics for other variables.

The calculation process for finding the median involves several internal steps: sorting the non-missing values in ascending order, counting the total number of non-missing observations, and then locating the middle value. If the count of observations is odd, the median is the single middle value. If the count is even, the median is the average of the two middle values. In the context of our 'points' column, the raw data points available are , which total 7 observations after excluding the missing value, resulting in a single middle value.

The following code snippet demonstrates the execution and resulting output when calculating the median of the 'points' column from our sample DataFrame:

```
# Find the median value of the points column using Series selection
```

```
df.median()
```

```
23.0
```

The calculated median value in the points column is indeed **23.0**. This indicates that 50% of the non-missing players scored 23 points or less, and 50% scored 23 points or more. This statistic provides a robust measure of central tendency for this metric, unaffected by the missing data point or any potential outliers.

## Handling Missing Values During Median Calculation

A crucial aspect of using the **median()** function in Pandas is its default behavior regarding missing data, often represented as Not a Number (NaN) or **pd.NA**. By default, the **median()** function is designed to ignore any missing or null values when calculating the final statistic. This means that the calculation only considers valid numerical entries, ensuring the resulting median is accurate based on the available data points.

This default handling is typically desirable in data analysis, as imputing or including null values in a percentile calculation can lead to skewed or misleading results. For instance, in our 'points' column, we had 8 potential observations, but one was missing. Pandas automatically adjusted the calculation base to the 7 available observations. If this default behavior were not implemented, the analyst would first need to manually clean the data by dropping or imputing missing values before calculating the median, adding an unnecessary step to the workflow.

However, analysts must be aware that this behavior can be controlled using the **skipna** parameter. By setting **skipna=False**, the function will return **NaN** if any missing values are present in the dataset, forcing the analyst to address data quality issues explicitly. For standard analysis, however, the default setting of **skipna=True** (which is implicit) is highly beneficial:

It ensures the calculation is based only on numeric column entries.

It prevents missing data from arbitrarily biasing the resulting median.

It simplifies the code required for robust statistical summaries.

## Example 2: Calculating the Median Across Multiple Columns

Often, data analysis requires assessing the central tendency across several related variables simultaneously. Pandas facilitates this by allowing the user to select multiple columns using a list of column names, effectively creating a subset DataFrame, and then applying the **median()** function. This method returns a Pandas Series where the index represents the column names and the values represent the median calculated for each corresponding column.

This simultaneous calculation is particularly useful for comparative statistics. For example, if we want to quickly compare the central performance metrics of 'points' and 'rebounds', calculating both medians in one operation provides a clean and efficient output. The underlying process remains the same as Example 1, but it is executed sequentially across all selected columns. The output format--a Series--is ideal for immediate interpretation or subsequent plotting.

In the following demonstration, we calculate the median for both the 'points' and 'rebounds' columns. Note that the median for 'points' is an integer (23.0), whereas the median for 'rebounds' is a float (8.5), which occurs because the 'rebounds' column has an even number of non-missing observations (8 total), requiring the average of the two central values:

```
# Find the median value of points and rebounds columns using list indexing  
df].median()
```

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

### Example 3: Finding the Median for All Numeric Columns

For a broad, initial statistical overview of a `DataFrame`, calculating the median for every numerical feature is highly efficient. Pandas makes this trivial: simply invoke the `median()` function directly on the entire `DataFrame` object without specifying any column names. Pandas intelligently identifies all valid `numeric columns` and automatically skips any columns containing strings, objects, or other incompatible data types.

This generalized approach is powerful for rapid exploratory data analysis (EDA). It provides a quick snapshot of the central tendency across the entire dataset, allowing the analyst to immediately compare the typical values of variables like 'points', 'assists', and 'rebounds' without needing to list them explicitly. If a `DataFrame` contained dozens of features, this method dramatically reduces the lines of code required for aggregation.

In our example `DataFrame`, 'points', 'assists', and 'rebounds' are all numerical. When `df.median()` is executed, it calculates the central value for each of these columns, resulting in the following comprehensive statistical summary:

**# Find the median value of all numeric columns in the DataFrame**

```
df.median()
```

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

The resulting output clearly shows that the median 'points' score is 23.0, the median number of 'assists' is 8.0, and the median number of 'rebounds' is 8.5. This allows for immediate comparative assessment of the player statistics.

### Advanced Use Case: Calculating Median with Grouping

While calculating the global median is valuable, complex data analysis often requires finding the median of a variable categorized by another, typically categorical, variable. This is achieved in Pandas using the `groupby()` method combined with the `median()` aggregation. The `groupby()` operation splits the `DataFrame` into groups based on unique values in a specified column, applies the median function to the desired numerical column within each group, and then combines the results into an output structure.

For instance, if our player data included a 'Team' column (e.g., 'Team A' and 'Team B'), grouping by 'Team' and then calculating the median 'points' would provide the central scoring tendency for

players on each respective team. This conditional aggregation is fundamental for segmenting analysis and understanding how statistical distributions vary across different subsets of the data. The resulting output is a Series or DataFrame indexed by the grouping variable.

The steps for performing this advanced calculation are intuitive and follow the standard Pandas aggregation pattern:

Use `df.groupby('grouping_column')` to define the groups.

Select the target numeric column (e.g., ).

Apply the aggregate function `.median()`.

This pattern ensures that the median is calculated locally for each subgroup, providing detailed, context-specific statistical measures, thereby moving beyond simple global summaries.

## Summary of Best Practices for Median Calculation

Calculating the median is a routine but critical step in data analysis, particularly when assessing data resilience to outliers. To maximize efficiency and ensure accuracy when using the `median()` function in Pandas, adherence to certain best practices is advised. These practices involve verifying data types, understanding the treatment of null values, and selecting the most appropriate method for the desired level of aggregation.

Key recommendations include:

**Data Type Verification:** Always confirm that the target columns are of a numeric data type (integer or float). Applying `median()` to object or categorical columns will typically raise an error or return NaN.

**Leverage Default Null Handling:** Rely on the default `skipna=True` setting unless the specific analytical requirement demands that the presence of null values should flag the entire calculation as invalid.

**Choosing the Scope:** Use `df.median()` for targeted single column analysis, `df].median()` for specific comparative analysis, and `df.median()` for quick, full-dataset summaries.

In conclusion, the `median()` function in Pandas is a cornerstone of robust statistical reporting. Its seamless integration with the DataFrame structure allows for immediate calculation of the middle value, providing an excellent, outlier-resistant measure of central tendency. Mastering these simple syntax patterns equips the data professional with the tools necessary to perform fast, reliable, and insightful exploratory data analysis using Python.