

# How to Order a PySpark DataFrame by Multiple Columns

Authored by  
**stats writer**

January 20, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Order a PySpark DataFrame by Multiple Columns*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=126667>

## The Necessity of Multi-Column Sorting in Large-Scale Data Analysis

When working with large datasets, efficient organization is paramount for effective analysis and reporting. In the realm of big data processing, especially when utilizing frameworks like [PySpark](#), data is typically distributed across clusters. Therefore, the sorting operation--which involves shuffling data--must be handled with careful consideration for performance. Sorting a [PySpark DataFrame](#) is a fundamental requirement, but often, sorting by a single field is insufficient to reveal the necessary hierarchical structure within the data.

For instance, consider a dataset tracking player statistics. We might first want to group results by **team**, then by **position** within each team, and finally by **points** scored to break ties among players sharing the same team and position. This requirement necessitates multi-column ordering, where [PySpark](#) processes the sort criteria sequentially. The framework must first stabilize the ordering based on the primary column; subsequent columns are then only used to refine the order among rows that share identical values in all preceding columns.

The mechanism [PySpark](#) uses for this complex operation is the powerful `orderBy` function, which is a method of the [DataFrame](#) API. Mastering the application of `orderBy` with multiple columns, including controlling the direction (ascending or descending) for each column, is a key skill for any data engineer or analyst working in the [Apache Spark](#) ecosystem.

### Understanding the PySpark `orderBy` Function

The `orderBy` transformation is essential for generating ordered results in distributed computing. While it is conceptually similar to the `ORDER BY` clause in traditional SQL, its execution within a distributed environment involves complex operations behind the scenes, such as data shuffling across nodes. When dealing with multi-column sorting, the order in which the columns are listed within the function's arguments dictates the priority of the sort. The column listed first serves as the primary sort key.

The function accepts either a single column name, a list of column names (strings), or a list of column expressions (`Column` objects) generated by functions like `asc()` or `desc()`. For simplicity and readability when applying consistent sort directions across all columns, providing a list of strings is the most common approach. However, as we will explore later, using column expressions is necessary when defining mixed sort orders (e.g., ascending on one column and descending on another).

By default, when you pass a list of column names, the [PySpark DataFrame](#) sorts all specified columns in **ascending** order. This behavior is intuitive for alphabetical and numerical data, moving from A to Z or smallest to largest, respectively. Understanding this default setting is critical to avoiding unexpected results when ordering complex datasets where consistency in sort direction is

not uniform.

## Basic Syntax for Ascending Order (The Default Behavior)

The most straightforward way to order a `PySpark DataFrame` by multiple columns relies on passing a Python list of column names directly to the `orderBy` function. Since the default behavior is ascending sort, no additional parameters are required for this standard implementation. The order within the list defines the hierarchy of sorting, ensuring a stable and deterministic result.

Below is the standard syntax for ordering by three distinct columns--`team`, `position`, and `points`--in their natural ascending order:

```
df.orderBy().show()
```

This execution chain first orders the rows entirely based on the values in the `team` column. Once the primary sort is complete, any rows sharing the same team value are then sorted according to the `position` column. Finally, if multiple players share both the same team and position, the `points` column acts as the tertiary sort key to definitively establish the final row order. This hierarchical structure ensures deterministic and precise sorting across the distributed data partitions.

## Controlling Sort Direction: Utilizing the `ascending` Parameter

While ascending order is often desired, many analytical tasks require sorting by the highest values first, necessitating a descending sort. To apply a uniform descending order across all specified sort columns, `PySpark` provides the optional `ascending` parameter within the `orderBy` function signature. This parameter accepts a boolean value (`True` for ascending, `False` for descending).

By setting the `ascending` parameter to `False`, you instruct the distributed sorting engine to reverse the default order for every column specified in the input list. This is extremely useful when, for example, ranking teams or players based on metrics where higher values are prioritized, and you want to see the best results first.

To achieve a multi-column sort where `team`, `position`, and `points` are all sorted in descending order (Z to A, or largest to smallest), the syntax is modified as follows. Note that the sort keys are still listed in order of priority, but the direction is universally reversed:

```
df.orderBy(ascending=False).show()
```

Crucially, when `ascending=False` is applied, it acts globally on the list of columns provided. The

`team` column will be sorted Z to A, and the `points` column will be sorted largest to smallest. This simplifies the syntax considerably when all columns require the same reversed sort direction without needing individual column expressions.

## Practical Demonstration: Setting Up the PySpark Environment

To illustrate these concepts effectively, we must first establish a working environment and create a reproducible sample dataset. This example uses a small dataset representing basketball players, featuring attributes like team, position, points scored, and assists. We utilize the `SparkSession` entry point, which is the gateway to using [Apache Spark](#) functionality through the Python API.

We define the structure of our data as a list of lists and specify the column names in a separate list. We then employ `spark.createDataFrame()` to transform this local Python structure into a distributed [PySpark DataFrame](#). Viewing the initial DataFrame confirms that the data is loaded correctly and is ready for the sorting transformations we intend to apply.

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.getOrCreate()
```

```
# Define sample basketball player data
```

```
data = ,
```

```
,
,
,
,
,
,
,
]
```

```
# Define column names
```

```
columns =
```

```
# Create the DataFrame
```

```
df = spark.createDataFrame(data, columns)
```

```
# View the initial, unsorted DataFrame
```

```
df.show()
```

```
+---+-----+-----+-----+
|team|position|points|assists|
+---+-----+-----+-----+
| A| Guard| 11| 4|
```

```
| A| Forward| 8| 5|
| B| Guard| 22| 6|
| A| Forward| 22| 7|
| C| Guard| 14| 12|
| A| Guard| 14| 8|
| B| Forward| 13| 9|
| B| Center| 7| 9|
+---+-----+-----+-----+
```

## Case Study: Sorting by Team, Position, and Points (Ascending)

We now apply the multi-column ascending sort using the simple list syntax. We prioritize sorting by **team**, followed by **position**, and finally by **points**. Since the `ascending` parameter is omitted, `PySpark` automatically defaults to ascending order for all three keys, ensuring a natural alphabetical and numerical progression.

The resulting DataFrame demonstrates the crucial hierarchical sorting principle. The primary sort organizes all 'A' teams together, followed by 'B', and then 'C'. Within Team 'A', the positions are sorted alphabetically ('Forward' before 'Guard'). Finally, within the 'A' / 'Forward' subgroup, the points are ordered from 8 to 22, showcasing the tie-breaking nature of subsequent sort keys.

### `df.orderBy().show()`

```
+---+-----+-----+-----+
|team|position|points|assists|
+---+-----+-----+
| A| Forward| 8| 5|
| A| Forward| 22| 7|
| A| Guard| 11| 4|
| A| Guard| 14| 8|
| B| Center| 7| 9|
| B| Forward| 13| 9|
| B| Guard| 22| 6|
| C| Guard| 14| 12|
+---+-----+-----+-----+
```

Analyzing the resultant order highlights the exact steps taken by the distributed engine:

First, the `team` column establishes the primary order (A to Z).

Second, for rows sharing the same team, the `position` column provides the secondary order

(alphabetically).

Third, for rows matching both team and position, the `points` column provides the final tertiary sort (from smallest to largest).

## Case Study: Applying Global Descending Order

Next, we will demonstrate how to reverse the entire sort hierarchy by applying the global `ascending=False` parameter. This transformation instructs `PySpark` to sort every specified column in reverse alphabetical or numerical order. This is a highly efficient way to retrieve results ordered from highest priority or latest entries downwards.

When applying `ascending=False`, the `team` column will be ordered Z to A (C, B, A), positions will be ordered Z to A (Guard, Forward, Center), and `points` will be ordered largest to smallest. This comprehensive reversal ensures that the highest-scoring players on the alphabetically last team will appear at the top of the resulting DataFrame structure.

`df.orderBy(ascending=False).show()`

```
+---+-----+-----+-----+
|team|position|points|assists|
+---+-----+-----+
| C| Guard| 14| 12|
| B| Guard| 22| 6|
| B| Forward| 13| 9|
| B| Center| 7| 9|
| A| Guard| 14| 8|
| A| Guard| 11| 4|
| A| Forward| 22| 7|
| A| Forward| 8| 5|
+---+-----+-----+

```

The resulting order confirms the successful reversal of the hierarchy:

The primary sort key, `team`, is now ordered from Z down to A.

The secondary key, `position`, is ordered in reverse alphabetical order (Guard is placed before Forward).

The tertiary key, `points`, is ordered from the largest value to the smallest value within identical team/position groups.

## Advanced Sorting: Using Column Expressions for Mixed Order

In more advanced scenarios, it is frequently necessary to mix sort directions within a single operation--for example, sorting a categorical column alphabetically (ascending) while sorting a numerical performance metric inversely (descending). Since the boolean `ascending` parameter applies globally to all columns in the list, we must instead provide the `orderBy` function with explicit column expressions using `asc()` and `desc()` functions.

These functions, often imported from `pyspark.sql.functions`, wrap the column name string, transforming it into a `Column` object that carries directional metadata. This technique bypasses the global `ascending` parameter entirely, providing fine-grained control over the sorting hierarchy. If we wanted to order by `team` ascending, but then rank within teams by `points` descending, the approach would look like this:

```
from pyspark.sql.functions import asc, desc
```

```
df.orderBy(asc('team'), desc('points')).show()
```

While this syntax is slightly more verbose than passing a simple list of strings, it is the standard and most robust method for handling complex, mixed-direction multi-column sorts in PySpark, ensuring that the final output aligns perfectly with intricate analytical requirements.

## Performance Considerations for Sorting Large DataFrames

While highly effective, the `orderBy` function is recognized as an expensive operation in terms of distributed computing resources. Sorting fundamentally requires a global understanding of the data, which necessitates a significant data shuffle--moving data partitions across the cluster nodes. This operation is often the primary bottleneck in large-scale [Apache Spark](#) applications.

When sorting by multiple columns, the complexity increases, but the main cost remains the global shuffle triggered by the initial sort key. To mitigate performance degradation, data engineers should minimize the number of sort operations. If a dataset is needed in multiple sorted views, consider persisting or caching the data after the initial sort, or leverage techniques like bucketing and partitioning during data storage that align with common sort keys to reduce future sorting overhead.

Always profile your application to ensure that the required sorting operations do not create an unacceptable bottleneck in the overall ETL or analytical pipeline. Efficient data partitioning, combined with judicious use of the `orderBy` function, is key to maintaining high throughput in a distributed environment.