

# How to Reorder Columns in a PySpark DataFrame

Authored by  
**stats writer**

January 1, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Reorder Columns in a PySpark DataFrame*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=110376>

In [PySpark](#), manipulating the structure of your data is a fundamental task, and reordering columns is often necessary for data governance, readability, or integration with other systems. The most straightforward and universally accepted technique for achieving column reordering involves the use of the powerful [select\(\) function](#), a core transformation available on the [PySpark DataFrame](#) API. This function allows users to project specific columns into a new DataFrame while defining the exact sequence in which they should appear.

For instance, if you have a DataFrame named `df` with columns "ColA", "ColB", "ColC", and you wish to change the order to "ColC", "ColA", "ColB", you would execute: `df.select("ColC", "ColA", "ColB").show()`. Understanding how to leverage `select()` is crucial, as it is not only used for reordering but also for projecting subsets of data and applying transformations simultaneously. This method ensures that the resulting DataFrame is clean and adheres precisely to the required output schema, which is vital when processing large volumes of [Big Data](#).

## Understanding Column Reordering in PySpark

Working with data at scale requires strict control over the schema and presentation. In [PySpark](#), column order, while technically not affecting computation results (as columns are referenced by name), is essential for final output presentation, integration with legacy systems, and user readability. When preparing data for machine learning models or writing data to external databases like specialized SQL systems or data warehouses, adhering to a predefined column sequence is often a strict requirement.

The primary mechanism for reshaping the [PySpark DataFrame](#) structure is the use of DataFrame transformations. Unlike simple column renaming or value manipulation, reordering involves defining a new projection that includes all desired columns listed in the exact sequence required. This operation is fundamentally achieved using `select()`, which is optimized by the underlying Apache Spark engine to be performed efficiently across distributed nodes.

We will explore two foundational methods for tackling column reordering: the manual definition of a specific sequence for targeted restructuring, and the programmatic approach for automatically sorting columns based on alphabetical order. Both techniques rely on the core [select\(\) function](#), but utilize different input parameters (explicit lists versus dynamically generated lists) to achieve their respective goals. Mastering these methods ensures control over data structure regardless of the DataFrame's size or complexity. You can use the following methods to reorder columns in a PySpark DataFrame:

### Method 1: Reorder Columns in Specific Order

```
df = df.select('col3', 'col2', 'col4', 'col1')
```

## Method 2: Reorder Columns Alphabetically

```
df = df.select(sorted(df.columns))
```

## Prerequisites Setup: Creating the PySpark DataFrame for Demonstration

Before diving into the reordering techniques, it is essential to establish a working example. We will define a sample PySpark DataFrame representing basketball team statistics. This example will clearly illustrate how column order shifts from the original definition to the transformed output, providing visual confirmation of the methods applied.

To initialize our environment, we must first import the necessary modules, specifically the `SparkSession`, which is the entry point for all Spark functionality when working in Python via PySpark. Once the session is created, we define the raw data and the corresponding schema names. This step is critical because the initial column names are what we will be referencing and rearranging later.

The following examples show how to use each method with the following PySpark DataFrame. The code block demonstrates the setup process, resulting in a DataFrame with the schema: **team**, **conference**, **points**, and **assists**. Note the inherent order of these columns as defined during the creation process, which is the starting point for all our subsequent reordering examples:

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

#define data
data = ,
,
,
,
,
]

#define column names
columns =

#create dataframe using data and column names
df = spark.createDataFrame(data, columns)

#view dataframe
df.show()
```

```

+----+-----+-----+-----+
|team|conference|points|assists|
+----+-----+-----+-----+
| A| East| 11| 4|
| A| East| 8| 9|
| A| East| 10| 3|
| B| West| 6| 12|
| B| West| 6| 4|
| C| East| 5| 2|
+----+-----+-----+-----+

```

### Example 1: Reordering Columns in a Specific, Custom Sequence

The most common requirement in data processing is to enforce a specific, manually defined order. This is achieved by passing an ordered list of column names directly to the `select()` function. The order in which the string literals are provided within the `select()` call dictates the final schema order of the resulting DataFrame.

This approach offers granular control, ensuring that critical identifier columns (like **team** or **conference**) appear first, followed by key metrics (like **points** and **assists**). It is especially useful when creating standardized output files or interfacing with external APIs that mandate a fixed schema structure for data ingestion. The syntax is extremely straightforward, relying purely on [Python](#) list ordering capabilities translated through the DataFrame API.

The core concept is that `select()` creates a new DataFrame based on the specified projection. It does not modify the original DataFrame in place, adhering to the immutability principle central to Apache Spark. To save the reordered structure, the result of the `select()` transformation must be assigned back to a variable, often the original DataFrame variable itself (e.g., `df = df.select(...)`), effectively replacing the older, unordered version with the newly structured data.

We can use the following syntax to reorder the columns in the DataFrame based on a specific order, moving **conference** to the front:

```

#reorder columns by specific order
df = df.select('conference', 'team', 'assists', 'points')

```

```

#view updated DataFrame
df.show()

```

```

+-----+-----+-----+-----+

```

```
|conference|team|assists|points|
+-----+-----+-----+-----+
| East| A| 4| 11|
| East| A| 9| 8|
| East| A| 3| 10|
| West| B| 12| 6|
| West| B| 4| 6|
| East| C| 2| 5|
+-----+-----+-----+-----+
```

The columns now appear in the exact order that we specified: **conference**, **team**, **assists**, and **points**.

## Handling Large DataFrames: Dynamic Reordering Techniques

While specifying a custom order is simple for DataFrames with few columns, it becomes cumbersome and error-prone when dealing with wide tables containing dozens or hundreds of attributes--a common scenario in Big Data environments. For these situations, dynamic reordering techniques that rely on list manipulation become indispensable, allowing developers to programmatically control the column sequence without manually listing every single column.

A typical dynamic scenario involves prioritizing a small set of key columns at the front while moving the remaining columns to the end in their original or alphabetical order. This can be achieved using basic Python list comprehension and operations. The first step is to identify the columns to prioritize. The second step is to calculate the remaining columns by subtracting the prioritized list from the full list of columns (available via `df.columns`).

For example, to place **team** and **conference** first, followed by all others in their original relative order, you would construct the final list: `prioritized_cols + remaining_cols`. This strategy provides flexibility, ensuring that even if new columns are added to the source data over time, the reordering script remains robust and does not require constant manual updates. This is crucial for maintaining automated ETL pipelines built using PySpark.

### Example 2: Reordering Columns Alphabetically

In situations where standardization is key, and the semantic order of columns is less important than their consistent placement, sorting columns alphabetically is a highly effective method. This is often utilized when preparing data for generic exploration tools or when consistency across multiple heterogeneous datasets is desired. Alphabetical reordering ensures that the schema is predictable regardless of the original data source or the order in which the columns were initially

loaded.

This method leverages [Python's](#) native `sorted()` function applied directly to the list of column names obtained from the DataFrame's `.columns` attribute. The `df.columns` property returns a standard Python list containing all column names as strings. Applying `sorted()` to this list generates a new list where the names are arranged in ascending alphabetical (lexicographical) order.

We can use the following syntax to reorder the columns in the DataFrame alphabetically:

```
#reorder columns alphabetically  
df = df.select(sorted(df.columns))
```

```
#view updated DataFrame  
df.show()
```

```
+-----+-----+-----+-----+  
|assists|conference|points|team|  
+-----+-----+-----+-----+  
| 4| East| 11| A|  
| 9| East| 8| A|  
| 3| East| 10| A|  
| 12| West| 6| B|  
| 4| West| 6| B|  
| 2| East| 5| C|  
+-----+-----+-----+-----+
```

The columns now appear in alphabetical order: **assists**, **conference**, **points**, and **team**.

## Performance Considerations and Best Practices

While column reordering using `select()` is a standard and safe transformation in PySpark, it is important to understand its performance implications, especially when working with massive datasets. The `select()` operation itself is a **narrow transformation**, meaning it does not require a full shuffle of the data across the cluster nodes. It operates primarily on the metadata (schema) and defines a new logical plan for the output DataFrame. The actual data movement only occurs when an action (like `.show()`, `.write()`, or `.collect()`) is called.

However, excessive use of transformations, even narrow ones, can increase the complexity of the underlying Catalyst Optimizer plan. For optimal performance, column reordering should ideally be performed as late as possible in the ETL pipeline, typically just before writing the final results to

storage. Reordering intermediate DataFrames unnecessarily adds overhead without contributing to data quality or core business logic.

Furthermore, when selecting columns, ensure that you only select the columns you actually need. If a DataFrame has 50 columns and you only need 10, reordering and selecting those 10 simultaneously is far more efficient than selecting all 50 and then dropping the unnecessary ones later. This principle of projection optimization is crucial in distributed computing frameworks like Spark, helping to minimize the amount of data processed and transferred throughout the cluster.

A brief summary of best practices for schema management in PySpark:

Always use the `select()` transformation for schema manipulation and column reordering; avoid deprecated or overly complex methods.

Minimize the number of times you reorder the same `DataFrame` throughout a single job execution. When using dynamic reordering techniques (like prioritizing certain columns), always verify that the resulting list includes every required column to prevent accidental data exclusion.

If you are working with extremely wide DataFrames (thousands of columns), consider techniques that abstract column selection, such as reading the schema from a template file rather than hardcoding lists, making the process easily manageable and scalable in Big Data environments.

## Advanced Scenario: Combining Transformations During Reordering

A significant advantage of using the `select()` function for column reordering is the ability to simultaneously perform data transformations, calculations, and column renaming within the same operation. Instead of merely listing column names, the `select()` function accepts expressions, allowing for powerful, consolidated data manipulation.

For example, if we wanted to reorder the columns while also renaming the `points` column to `total_score`, or applying a simple calculation to the `assists` column, we could do this within the same select statement. We use `F.col()` or `df` combined with the `.alias()` method for renaming, or simply algebraic operations for calculations.

This method significantly streamlines the code, eliminating the need for multiple sequential transformation steps (e.g., using `withColumnRenamed` followed by `select` for reordering). The transformation and reordering are executed logically together, resulting in a cleaner execution plan. If we wanted to keep our custom order (conference, team, assists, points) while converting points to scores and adding a new calculated column called `ratio` (points/assists), the select statement would look like this:

```
from pyspark.sql import functions as F
```

```
df_transformed = df.select(
```

```
F.col("conference"),
F.col("team"),
F.col("assists"),
F.col("points").alias("total_score"), # Renaming points during selection
(F.col("points") / F.col("assists")).alias("ratio") # Creating a new column
)
df_transformed.show()
```

This demonstrates the utility of `select()` beyond simple schema restructuring. By integrating transformations, we maintain control over the final column sequence while ensuring all necessary data preparation steps are executed efficiently, reflecting a powerful pattern for structured data manipulation in [PySpark](#).

## Conclusion and Further Learning

Reordering columns in a [PySpark DataFrame](#) is a foundational task easily accomplished using the `select()` transformation. Whether your requirement is to enforce a specific, manually defined sequence or to apply an automated alphabetical sorting, the solution involves providing `select()` with an ordered list of column identifiers. We have seen that this process is highly flexible, allowing for dynamic list generation crucial for managing large-scale data schemas effectively.

The ability to dynamically generate column lists using Python's core functionalities--such as list comprehensions or the `sorted()` function--allows developers to write concise and maintainable code that scales well to complex [Big Data](#) challenges. Always prioritize code clarity and efficiency by performing reordering only when necessary and by combining it with other necessary transformations where logical.

To further enhance your skills in PySpark and data manipulation, we encourage exploring related topics concerning data type casting, column dropping, and complex data restructuring techniques. Understanding these advanced DataFrame operations is key to becoming proficient in Spark development.

The following tutorials explain how to perform other common tasks in PySpark: