

# Learning How to Rename Columns in PySpark DataFrames: A Step-by-Step Guide

Authored by  
**Mohammed loot**

November 11, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Learning How to Rename Columns in PySpark DataFrames: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16648>

## Introduction to Column Renaming in PySpark

When working with large-scale data processing using [Apache Spark](#), specifically through its [Python](#) API, [PySpark DataFrame](#) manipulation is a daily necessity. Renaming columns is a fundamental operation required for data standardization, improving readability, integrating datasets with differing naming conventions, or preparing data for machine learning models. Fortunately, PySpark provides several straightforward and efficient methods to handle column renaming, whether you need to change a single column, a specific subset, or overhaul the entire schema. The choice of method often depends on the scale of the renaming task and the desired level of programmatic control.

The primary technique for renaming individual columns involves the built-in function [withColumnRenamed](#). This method is highly favored due to its simplicity and clarity when dealing with minor modifications. For scenarios requiring a complete schema replacement, the [toDF](#) function offers a powerful, albeit more aggressive, approach that demands careful adherence to column order. Understanding these core techniques is essential for effective data wrangling within the Spark ecosystem.

We will explore three distinct methods for renaming columns in a [PySpark DataFrame](#), providing practical examples for each approach. These methods are designed to cover nearly every renaming requirement you might encounter during data preparation:

**Method 1: Rename One Column** using the intuitive [withColumnRenamed](#) function.

**Method 2: Rename Multiple Columns** by chaining calls to [withColumnRenamed](#).

**Method 3: Rename All Columns** using the [toDF](#) function, which is ideal for applying a new, complete schema.

## Prerequisites: Setting up the PySpark DataFrame

To demonstrate these renaming techniques effectively, we first need to establish a base [PySpark DataFrame](#). This DataFrame represents typical structured data containing information about sports teams, including their conference, points scored, and assists. The setup involves initializing a `SparkSession`, defining the raw data, and specifying the initial column names, which will be the target of our renaming operations throughout the subsequent examples.

The initial column names--`team`, `conference`, `points`, and `assists`--are relatively clear, but in a real-world scenario, you might want to adjust them for consistency (e.g., changing `conference` to `conf` for brevity or `team` to `team_id` for specificity). The following code block illustrates the standard process of creating and displaying the starting DataFrame structure.

```
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|
+---+-----+-----+-----+
```

This initial structure serves as the foundation for all subsequent examples. PySpark transformations, including renaming, operate under the principle of immutability, meaning that each renaming operation returns a new DataFrame rather than modifying the original one in place. This is a critical concept for maintaining data integrity and traceability in large-scale distributed computing environments like [Apache Spark](#).

## Method 1: Renaming a Single Column (withColumnRenamed)

The simplest and most recommended approach for renaming an isolated column is utilizing the [withColumnRenamed](#) function. This function takes exactly two arguments: the name of the existing column and the desired new name. It is highly readable and minimizes the risk of unintended

changes to the schema, as it only focuses on the specified column.

For example, if the column `conference` is too verbose, we can easily shorten it to `conf`. This transformation is applied directly to the DataFrame object and, as is standard in PySpark, returns a new DataFrame object with the updated schema. It is important to assign the result back to a variable (usually the original DataFrame variable, `df`) to capture the change.

```
#rename 'conference' column to 'conf'
```

```
df = df.withColumnRenamed('conference', 'conf')
```

```
#view updated DataFrame
```

```
df.show()
```

```
+----+----+-----+-----+
|team|conf|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|
+----+----+-----+-----+
```

As demonstrated by the output, only the `conference` column has been successfully renamed to `conf`. This method is idempotent for individual column changes and provides the cleanest syntax when only minor schema adjustments are necessary across a [PySpark DataFrame](#).

## Method 2: Renaming Multiple Columns (Chaining Transformations)

When the requirement extends beyond a single column to a specific subset of columns, the most efficient way to use `withColumnRenamed` is through method chaining. Since each call to the function returns a new DataFrame, we can sequentially apply multiple renaming operations in a single, fluent expression. This approach is highly favored in the [Python](#) community for its readability and conciseness, avoiding the need for intermediate variable assignments.

For instance, we might want to standardize both the `conference` column (to `conf`) and the `team` column (to `team_name`). We simply chain the two `withColumnRenamed` calls together, ensuring that the second operation acts upon the result of the first. It is good practice to use indentation, as shown below, to make the chained operations easily distinguishable within the code block.

```
#rename 'conference' and 'team' columns
df = df.withColumnRenamed('conference', 'conf')
.withColumnRenamed('team', 'team_name')
```

```
#view updated DataFrame
```

```
df.show()
```

```
+-----+----+-----+-----+
|team_name|conf|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|
+-----+----+-----+-----+
```

This technique confirms that both specified columns, `team` and `conference`, have been successfully updated to `team_name` and `conf`, respectively, while the remaining columns (`points` and `assists`) are left untouched. Alternatively, for renaming a large number of columns programmatically, one could define a mapping dictionary and iterate through it, applying `withColumnRenamed` within a loop, although chaining remains the cleanest solution for a fixed, small set of changes.

### Method 3: Renaming All Columns (toDF)

In scenarios where the entire schema needs a complete overhaul, perhaps to adhere to a strict organizational standard or integrate with another system, providing a full list of new column names is more efficient than chaining dozens of individual renaming calls. This is where the `toDF` function proves invaluable. The `toDF` function is specifically designed to assign a new sequence of column names to the existing data structure.

A critical aspect of using `toDF` is the requirement that the list of new column names must exactly match the number of columns currently present in the DataFrame. Furthermore, the order in which the new names are supplied directly maps to the existing column order. Failure to provide the correct count or order will result in a DataFrame where data fields are incorrectly labeled, potentially leading to serious downstream data integrity issues.

In this example, we define a list called `col_names` containing four elements, which corresponds

precisely to the four columns in our original DataFrame. We then use the Python asterisk operator (\*) to unpack this list directly into the `toDF` function call, replacing the entire schema instantly.

### #specify new column names to use

```
col_names =
```

```
#rename all column names with new names
```

```
df = df.toDF(*col_names)
```

```
#view updated DataFrame
```

```
df.show()
```

```
+-----+-----+-----+-----+
|the_team|the_conf|points_scored|total_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|
+-----+-----+-----+-----+
```

The output confirms that every column name has been updated according to the list provided. While powerful, `toDF` is generally reserved for comprehensive schema replacements where the risk of misordering is minimized, perhaps immediately after data ingestion or a sequence of transformations that guarantee a consistent column order.

## Summary and Best Practices for PySpark Renaming

Selecting the appropriate method for column renaming in [Apache Spark](#) depends largely on the scope of the required changes. For targeted modifications, `withColumnRenamed` is the uncontested champion, offering safety, clarity, and ease of use. It is the default choice for improving data quality one column at a time or even for small batches via chaining. It ensures that only the intended columns are affected, maintaining the integrity of the rest of the [PySpark DataFrame](#) structure.

Conversely, `toDF` is a powerful utility best suited for complete schema definition tasks. While it requires strict attention to the mapping order, it significantly reduces boilerplate code when dealing with dozens or hundreds of columns that need new, standardized names. A common best practice is to extract the current column order (using `df.columns`) before defining the new list of names to ensure accurate mapping, thereby mitigating potential errors introduced by relying solely on visual

inspection.

Mastering these renaming techniques is crucial for efficient data preparation in PySpark. Regardless of the method chosen, always remember the immutable nature of Spark DataFrames--every operation creates a new state, which must be captured by reassignment. For further exploration into common data manipulation tasks, the following resources provide excellent documentation and examples for extending your PySpark expertise.

## **Additional Resources**

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